

Tag-based scheduling system for digital communication switch

Publication number: US5455825

Publication date: 1995-10-03

Inventor: LAUER HUGH (US); SHEN CHIA (US); GHOSH ABHIJIT (US)

Applicant: MITSUBISHI ELECTRIC RESEARCH L (US)

Classification:

- international: **H04Q3/00; H04L12/56; H04Q3/52; H04Q3/00; H04L12/56; H04Q3/52;** (IPC1-7): H04L12/56

- European: H04L12/56E3

Application number: US19940234385 19940428

Priority number(s): US19940234385 19940428

Also published as:

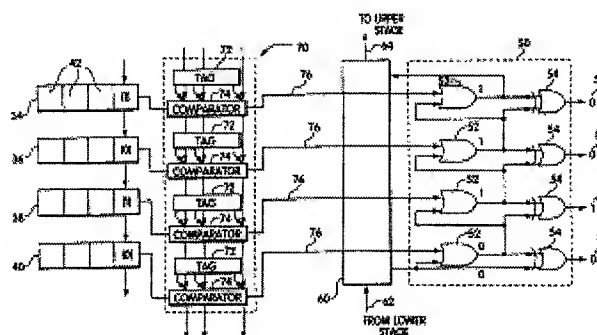


JP8056230 (A)

[Report a data error here](#)

Abstract of **US5455825**

A switch for digital communication networks includes a queuing system capable of implementing a broad class of scheduling algorithms for many different applications and purposes, with the queuing system including means for providing numerical tags to incoming cells or packets, the values of the tags being calculated when incoming cells or packets arrive at the switch. A queue and search module is provided to select cells or packets for transmission based on these tags. The combination of the tags and the queue and search module enables simple and fast implementations of a wide variety of scheduling algorithms, including algorithms for supporting communication traffic with real time requirements, continuous media such as audio and video, and traffic requiring very fast response. Furthermore, multiple classes of traffic are supported in a single network switch, each class having its own scheduling algorithm and policy. The queue and search module is designed for VLSI implementation, and in one embodiment supports an ATM switch with 16 ports, each port operating at 622 megabits per second.



Data supplied from the **esp@cenet** database - Worldwide

(19)日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11)特許出願公開番号

特開平8-56230

(43)公開日 平成8年(1996)2月27日

(51)Int.Cl. ⁶	識別記号	片内整理番号	F I	技術表示箇所
H 0 4 L 12/28				
12/56				
H 0 4 Q 3/00				
	9466-5K	H 0 4 L 11/ 20	H	
	9466-5K		G	
審査請求 未請求 請求項の数22 O L (全 22 頁) 最終頁に続く				

(21)出願番号 特願平7-103908

(22)出願日 平成7年(1995)4月27日

(31)優先権主張番号 0 8 / 2 3 4 3 8 5

(32)優先日 1994年4月28日

(33)優先権主張国 米国 (U S)

(71)出願人 000006013

三菱電機株式会社

東京都千代田区丸の内二丁目2番3号

(72)発明者 フー ラウア

アメリカ合衆国 マサチューセッツ州 コ
ンコルド ブローダーロード 69

(72)発明者 チア シェン

アメリカ合衆国 マサチューセッツ州 ソ
マービル モリソンアベニュー 169

(72)発明者 アブヒジット ゴーシュ

アメリカ合衆国 カリフォルニア州 パー
クレイ スラターレーン 42

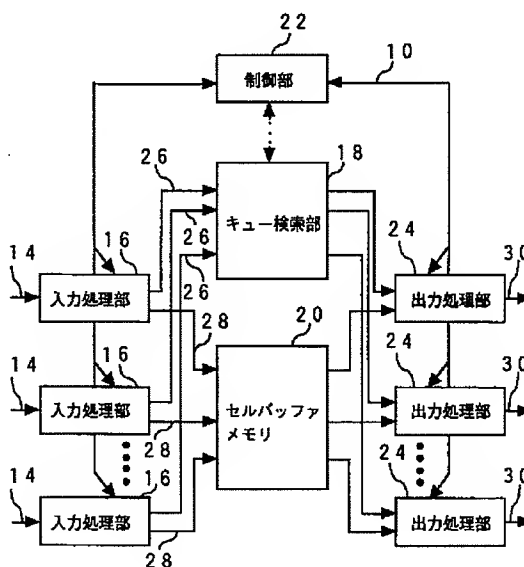
(74)代理人 弁理士 高田 守 (外4名)

(54)【発明の名称】 スイッチングシステム

(57)【要約】

【目的】 ダイナミック・スケジューリング・アルゴリズムをサポートできるタグベースのスイッチングシステムを得る。

【構成】 入力したセルに宛先ビットと優先度を示すタグを付加し、タグの値により出力順を決定する。タグの値の比較は、木構造で配置された比較器により並列に行うので、処理速度の速いスイッチングシステムが実現できる。



1

【特許請求の範囲】

【請求項1】 複数の入力リンクと複数の出力リンクを持ち、各入力リンクによりネットワーク内の1以上の宛先を示す宛先情報を含むヘッダフィールドとデータとを格納したセルを受信するスイッチと、

上記ヘッダフィールドの宛先情報を、各ビットが各出力リンクに対応しておりそのビットが有意状態である時そのセルが最終宛先に到着するようにそのビットに対応する出力リンクに対してそのセルを出力することを示している宛先ビットのベクトルに変換する変換手段と、

各セルに対して、上記宛先ビットのベクトルを付加するベクトル付加手段と、

上記ベクトルを付加されたセルを1つ以上の出力リンクに接続する接続手段とを備え、

上記接続手段は、上記ベクトルを付加されたセルを到着順に格納するキューイング手段と、キューイング手段に格納されたセルの中から到着順とは異なる順に各出力リンクに対して各出力リンクに対応する宛先ビットが有意状態のセルを選択する選択手段とを備え、

上記選択手段は、セルの出力をスケジュールするスケジュール手段を有し、上記スケジュール手段は、選択順を決定するアルゴリズムを記憶する手段と、

スイッチの動作に先だってそのアルゴリズムを設定する手段と、

スイッチの動作中にそのアルゴリズムを再設定する手段と、

上記各出力リンクに対してセル送信時間内に上記アルゴリズムを実行する手段を有しており、

上記スイッチは、更に、上記出力リンクに選択したセルを出力する出力手段を備えたことを特徴とするデジタル通信ネットワーク用のスイッチングシステム。

【請求項2】 上記出力手段は、上記ベクトルが除去されたセルを出力することを特徴とする請求項1記載のスイッチングシステム。

【請求項3】 上記アルゴリズムは、ダイナミック・プライオリティ・スケジューリング・アルゴリズムと、スタティック・プライオリティ・スケジューリング・アルゴリズムと、ラウンド・ロビン・アルゴリズムと、これらの組み合わせのいずれかから選択されることを特徴とする請求項1記載のスイッチングシステム。

【請求項4】 上記スケジュール手段は、各出力リンクにセルを出力する順番を決定するために設定されたスケジュールリングアルゴリズムに基づいて、キューイング手段の中に格納されたセルの出力順序を示す値をタグ値として計算し、このタグ値を有するタグをセルに付加するタグ付加手段を有し、上記選択手段は、タグ値と宛先ビットの値に基づいて、キューイング手段の中に格納されたセルをサーチするサーチ手段とを備え、選択手段は、各出力リンクに対応する宛先ビットの値が有意状態であるセルの中で最も小さいタグ値を持つセルをサーチして

2

選択し、もし、最も小さいタグ値を持つセルが2以上存在する場合に、先に到着したセルを選択することを特徴とする請求項1記載のスイッチングシステム。

【請求項5】 上記キューイング手段は、先頭から末尾に至る複数のレジスタからなるFIFO回路を備え、各レジスタは、出力を待っているセルを示すものであり、各レジスタは、上記宛先ビットのベクトルとタグ値のバイナリーコードを保持し、各レジスタは、先頭から末尾の方向でセルの到着順を示していることを特徴とする請求項4記載のスイッチングシステム。

【請求項6】 上記サーチ手段は、上記レジスタの数と同じ数の複数の比較回路を備え、各比較回路は、各レジスタに対応して設けられ、各比較回路は、対応したレジスタのタグ値と末尾側にある次のレジスタに対応した比較回路からの出力値を比較するようにリニアに配置され、各比較回路は、対応するレジスタのタグ値が末尾側にある次に比較回路から出力されたタグ値より小さい場合であって、かつ、選択した出力リンクに対する宛先ビットが有意状態である場合に、出力ビットをオンにすることにより、上記比較回路は、選択した出力リンクに対して最小のタグ値を伝搬させ、待ち行列の中からセルを選択することを特徴とする請求項5記載のスイッチングシステム。

【請求項7】 上記サーチ手段は、上記レジスタの数の半分に当たる複数の比較回路と追加の比較回路を備え、上記比較回路を階層的に配置し、ルートと枝と葉を持つツリーを構成し、葉と枝のサブセットによりサブツリーを構成し、葉に配置された各比較回路は、2つの隣り合うレジスタのタグ値を比較して小さいタグ値を有するレジスタを識別し、上記追加の比較回路は、2つのサブツリーからのタグ値を比較し、選択した出力リンクに対して、最小のタグ値を持つレジスタが存在しているサブツリーを識別し、上記ツリーは、選択した出力リンクに対して、最小のタグ値を伝搬することを特徴とする請求項5記載のスイッチングシステム。

【請求項8】 複数の出力リンクと出力リンクを示すアドレスを含んだセルを受信する複数の入力リンクと、上記アドレスを宛先ビットに変換する変換手段と、上記入力リンクのセルを少なくとも1つの出力リンクに接続する接続手段とを備え、

上記接続手段は、入力リンクから出力リンクへセルの出力をスケジュールするスケジュール手段を備え、スケジュール手段は、

各出力リンクにセルを出力する順番を決定するために、予め設定されたスケジュールリングアルゴリズムに基づいて、どのセルを出力すべきかを示す値を生成し、生成した値をタグ値として有するタグをセルに付加するタグ手段と、

共通キューと、

上記共通キューに、各セルに対して宛先ビットとタグ値

を挿入する手段と、

上記宛先ビットとタグ値に基づいて、上記共通キューから、対応している宛先の中から最小のタグ値を持つセルをサーチするサーチ手段と、

宛先に対応した出力リンクに対して最小のタグ値を持つセルを出力する出力手段とを備えたことを特徴とするデジタル通信ネットワーク用のスイッチングシステム。

【請求項9】 上記サーチ手段は、最小のタグ値をリアルに評価していくリニアサーチ手段を有していることを特徴とする請求項8記載のスイッチングシステム。

【請求項10】 上記サーチ手段は、木構造を用いて同時に複数対のセルを評価していくログリズミックサーチ手段を有していることを特徴とする請求項8記載のスイッチングシステム。

【請求項11】 上記サーチ手段は、リニアサーチ手段とログリズミックサーチ手段を結合したサーチ手段を有していることを特徴とする請求項8記載のスイッチングシステム。

【請求項12】 上記タグ手段は、先頭から末尾までのキューを形成するように配置された一連のタグレジスタと、

複数の比較器と、

上記タグレジスタと比較器を接続する接続手段と、

上記セルに対応するタグレジスタにより識別できるように記憶する記憶手段と、

記憶されたセルのタグ値をタグレジスタに設定する設定手段とを備え、

上記比較器は、タグレジスタのタグ値と1つ前の比較器からの出力値を比較して小さい方の値を決定するとともに、

小さい方の値と宛先ビットを出力する出力手段と、

小さい方の値と宛先ビットに基づいて、記憶されたセルを識別する識別手段とを備えたことを特徴とする請求項8記載のスイッチングシステム。

【請求項13】 上記比較器は、一連のビットから構成された出力を持ち、この一連のビットの中に、宛先ビットが有意状態であり、同じ宛先ビットが有意状態であるキュー中の他のセルのタグ値よりも小さいタグ値を持つことを示す特定ビットを有しており、

各比較器は、各タグ値とキューの末尾方向にあるタグの最小値とを比較することにより、宛先別に各セルのプライオリティを設定することを特徴とする請求項12記載のスイッチングシステム。

【請求項14】 上記比較器は、2つのタグ値を比較する手段と、小さい方の値を出力する手段とを有し、1つのタグ値に対して1つの追加の入力ビットを備え、上記比較器の出力は、比較対象となるタグ値を示す追加の入力ビットにより条件付けられていることを特徴とする請求項12記載のスイッチングシステム。

【請求項15】 上記追加の入力ビットは、宛先ビット

が有意状態になっているセルのタグ値を特定するものであることを特徴とする請求項14記載のスイッチングシステム。

【請求項16】 上記スイッチングシステムは、更に、最小のタグ値を持つセルが2以上存在する場合に、キューの先頭に近いセルを選択する手段を備えたことを特徴とする請求項12記載のスイッチングシステム。

【請求項17】 上記タグ手段は、枝を経由してルートに連結された葉を持った階層の木構造を有し、各枝は、一对のタグレジスタと隣り合うタグレジスタの間に設けられたタグレジスタ用比較器と、2つのタグレジスタ用比較器からの出力を入力する追加の比較器を有しており、上記追加の比較器は、宛先ビットが有意状態になっており、小さい方のタグ値を持つタグレジスタに対応するセルを識別する手段を有していることを特徴とする請求項8記載のスイッチングシステム。

【請求項18】 上記タグレジスタ用比較器は、小さい方のタグ値と2つの比較ビットを出力し、2つの比較ビットの内、一方の比較ビットは一方のセルが小さい方のタグ値を持っていることを示し、2つの比較ビットの内、他方の比較ビットは他方のセルが小さい方のタグ値を持っていることを示すものであり、上記選択手段は、小さい方のタグ値を持っているセルを識別するために、宛先ビットと比較ビットを入力する複数のANDゲートを有し、複数のANDゲートは、ルートにおいてただ1つのANDゲートが有意な出力を有するように構成されており、このANDゲートの有意な出力により出力リンクに出力されるセルを示すことを特徴とする請求項17記載のスイッチングシステム。

【請求項19】 上記ANDゲートは、そのANDゲートに対応したセルが適切な宛先ビットを有しているかを判定するために用いられるものであり、上記タグレジスタ用比較器は、異なるANDゲートへ接続された出力ビットを有し、このタグレジスタ用比較器からの出力ビットは、対応するセルの宛先ビットとANDを取られ、その結果の出力が有意状態である場合は、対応するセルが宛先ビットが有意状態であり、最小のタグ値を持つセルであることを示すことを特徴とする請求項18記載のスイッチングシステム。

【請求項20】 上記スイッチングシステムは、更に、2以上のセルが最小のタグ値を有する場合、先着のセルを識別する手段を有することを特徴とする請求項19記載のスイッチングシステム。

【請求項21】 上記タグレジスタ用比較器は、先頭と末尾を持つキューを定義し、上記追加の比較器は、2入力から小さい値を出力するとともに、キューの先頭方向にある2入力の値がキューの末尾方向にある2入力の値以下である場合に、有意状態を出力することを特徴とする請求項17記載のスイッチングシステム。

【請求項22】 上記スイッチングシステムは、更に、

5

上記木構造のルート方向への上位層レベルにおいて、下位層レベルのどの枝が小さい方のタグ値を有しているか決定することにより小さい方のタグ値を有する枝を識別する識別手段を有し、上記識別手段は、各タグに対応してそのタグの宛先ビットとそのタグが最小のタグ値を有することを示すビットとのANDを取るためのANDゲートを有していることを特徴とする請求項17記載のスイッチングシステム。

【発明の詳細な説明】

【0001】

【産業上の利用分野】この発明は、セルやパケットがノード間でネットワークを転送されるデジタル通信ネットワークに関するものである。さらに、この発明は、ネットワーク中でセルが、あるノードから次のノードへ転送される順序をスケジューリングするための改良された方式を備えたデジタル通信ネットワークに関する。

【0002】

【従来の技術】一般に、デジタル通信ネットワークでは、メッセージや情報の流れがパケットやセルと呼ばれる連続した小単位に細分化され、このセルやパケットが、ノードからノードへ伝送される。また、各ノードにおいて、スイッチがセルやパケットを伝送する順序と伝送する次のノードとを選択する。その結果、デジタル情報はその最終宛先へ適時に到着する。これらのネットワークスイッチは、様々な特性のネットワーク通信をサポートできるのが望ましい。例えば、適時な到着をリアルタイムに厳しく保証することを要求する通信、音声やビデオ用の連続的なメディア通信、早急な返答を要求する通信が含まれる。

【0003】重要なデジタル通信ネットワークの1つとして、ATM (Asynchronous Transfer Mode: 非同期伝送モード) ネットワークがある。ATMネットワークは、ネットワーク中のあるポイント (ノード) から、あるほかのポイント (ノード) へデータを伝送するのに用いられる。ここで、データや情報は、連続した固定サイズの小さなセルに細分化されて、ネットワークのノード間で伝送される。ここで述べているノードは、ネットワークの複数のノード間でパケットやセルを高速にスイッチングしたり、ルーティングしたりするATMスイッチを含む。ATMネットワークの一般的な原理は、J・ブライアン・リールズとダニエル・C・スウィンハートによる、「発生するギガビット環境とローカルATMの役割」(IEEEコミュニケーションマガジン、30巻、4類、1992年4月、52-58頁)、及びC・ラムによる、「ATM方式への高速化」(ユニックス レビュー、10巻、10類、1992年10月、29-36頁) に述べられている。

【0004】ATMネットワークの重要な課題の1つは、各スイッチによって伝送されるセルを、スケジューリングすることである。セルは通常、各スイッチの待ち

6

行列内に、バッファリングされる。セルの渋滞 (輻輳) がないものとする、これらのセルは、スイッチで入力リンクから受信され、直ちに出力リンクをとおして別の宛先へ送信される。しかし、セルが複数の入力リンクを通して到着し、同じ出力リンクへ同時に送信されなければならない場合、セルが望ましい順序で伝送されるには、セルの待ち行列 (キュー) をつくる必要がある。

【0005】どのセルがどの時間にどの順序で伝送されるかのスケジューリングを調整するためには、FIFO (First In First Out: 先入れ先出し) 方式による順序付けを使うのが普通である。FIFO方式においては、スイッチに、先に到着するセルが、次に到着するセルに優先して伝送される。リアルタイムアプリケーションをサポートするネットワークの場合、セルに優先度がつけられ、セルの優先度によって、セルは別々の待ち行列に記憶される。その後、セルは、その別個の待ち行列の優先度に従って、指示された順序で伝送される。これらの単純な方式は、限られた数と種類のリアルタイムアプリケーションのみをサポートしうる。なぜならば、FIFO方式と優先度スケジューリング機能は、データ損失のない適時な伝送に対して、限られた保証のみを提供できるからである。例えば、ATMローカル・エリア・ネットワーク、あるいは、ATMワイド・エリア・ネットワークの初期世代のスイッチは、単に、入力ネットワークリンクから出力ネットワークリンクへセルを転送するというスケジューリングをしただけで、大変簡易なスケジューリングアルゴリズムを備えている。通信は、ATM法則によるFIFO順序付けで処理されるが、少数の、通常2種類の優先順位が、リアルタイム通信要求を持つアプリケーションをサポートするために、かなり限定された方法で提供されている。静的に割り当てられた優先度は、近年の進んだ変化に富む多種多様なアプリケーションを持つ中位サイズのローカル・エリア・ネットワークに対してさえも、ほとんど十分でない。さらに、アプリケーションが音声やビデオなど連続的なメディアを伝送する必要がある場合、もしくは、アプリケーションが実時間で期待した返答を必要とする場合、ネットワークサービスの質や適時性に関して予測や保証をすることは実質的に不可能である。

【0006】過去20年にわたって、実時間やマルチメディア・コンピューティング及び通信に関して十分な研究がなされてきた。また、多くのスケジューリングアルゴリズムが細部にわたり発明され、研究されてきた。例えば、J・ジャンゴク・パーとT・スーダによる「ATMネットワークにおける通信量制御スキームとプロトコルの分析」(IEEE会報79巻、第2版、1991年2月、170-189頁) は、ATMスイッチにおいて一般的なスケジューリング機能によって、実現されな

ればならない多数の通信コントロールスキームを説明している。さらに、H・ジャングらは、「転送速度に基づくサービス原則の比較」(H・ジャングとS・ケシャブ、ACM SIGCOMM会報、'91、チューリッヒ、1991年9月)、「レート制御静的優先度待ち」(H・ジャングとD・フェラーリ、国際コンピュータサイエンス学会技術報告#TR-92-003、カリフォルニア、バークレー)の論文で、転送速度に基づくサービス原則と転送速度により制御された静的優先度キューイングを説明している。スケジューリングアルゴリズムのさらなる例は、以下の論文、論説に記載されている。W・A・ホーン、「簡易なスケジューリングアルゴリズム」(ナール リサーチ ロジスティックオータリ、21巻、1974年、177-185頁)、J・R・ジャクソン「最大遅延量を最小にするプロダクションラインのスケジューリング」(マネージメント サイエンス リサーチ プロジェクト、調査報告43、UCLA、1955年1月)、C・R・カルマネク、H・カナキア、S・ケシャブ、「高速ネットワークのためのレート制御サーバー」(IEEEグローバル テレコミュニケーション コンファレンス、カリフォルニア、サンディエゴ、1990年12月、300.3.1-300.3.9頁)、A・クマー、J・バレク「統合サービスネットワークにおけるフロー制御への汎用プロセッサのシェアリングアプローチ」(pHD論文、MIT、1992年2月)、H・T・カンク、A・チャップマン、「ATMネットワークに対するFCVC(Flow-Controlled Virtual Channels:フロー制御バーチャルチャネル)提案」(ネットワークプロトコルに関する1993年国際会議会報、カリフォルニア、サンフランシスコ、1993年、10月19-22日)、C・L・リュウとJ・W・レイランド、「ハードの実時間環境におけるマルチプログラミングのためのスケジューリングアルゴリズム」(ACMジャーナル、20巻、第1版、1973年)、L・ジャング、「バーチャルクロック:パケットスイッチネットワークのための新しい通信量コントロールアルゴリズム」(ACM SIGCOMM会報、ペンシルバニア、フィラデルフィア、1990年9月、19-29頁)、L・ジャング、「バーチャルクロック:パケットスイッチネットワークのための新しい通信量コントロールアルゴリズム」(コンピュータシステムにおけるACMトランザクション、9巻、第2版、1991年5月、101-124頁)、Q・ジェンク、K・シン、「2地点間パケットスイッチネットワークにおける実時間チャネルの確立能力について」(コミュニケーションにおけるIEEEトランザクション、1994年3月)。

【0007】これらのアルゴリズムの多くは、ATMネットワークにおける実行には、不適当とされてきたことに注意しなければならない。なぜならば、こうしたアル

ゴリズムは、一般的な環境には、専門的すぎるからである。また、キューのサーチに時間がかかりすぎるからである。

【0008】しかし、ATMスケジューリングのためのVLSIシーケンサーチップが、H・ジョナサン・カオとネクデット・アズンによる論文、「ATM通信シェイパーと行列管理のためのVLSIシーケンサーチップ」(IEEEソリッドステイトサーキットジャーナル、27巻、第11版、1992年11月)において述べられている。このシステムにおいて、セルは出力ポートにスイッチングされる。各ポートは、それ自体の待ち行列とシーケンサーチップを含んでいる。各ポートにおいて、セルは、シーケンサーチップによって優先順位に分類される。このシステムは、改良された優先度スケジューリングを実現しているが、スイッチの各出力ポートに対して別々の分類回路を要求する。その結果、コストが増加し、スイッチのデザインの柔軟性が減少する。また、性能低下の可能性もある。

【0009】共用バッファスケジューリング方式を用いたATMスイッチに関する重要な提案が、H・近藤ほか、K・大島ほかによる以下の2つの論文で述べられている。H・近藤、H・山中、M・石脇、Y・松田、M・中谷、「ATMスイッチLSIのための効率的なセルフタイムの行列アーキテクチャー」(カスタム インテグレートド サーキット コンファレンス、サンディエゴ、1994年5月)、K・大島、H・山中、H・斉藤、H・山田、S・小浜、H・近藤、Y・松田、「STSタイプの共用バッファリングとそのLSI実現に基づく新しいATMスイッチアーキテクチャー」(国際スイッチングシンポジウム'92会報、日本、横浜、1992年10月、359-363頁)。このATMスイッチデザインの重要な要素は、共通の待ち合わせ機能と、統計的な多重化を進展させるため共用バッファと、より迅速な性能とより低いコストである。このスイッチは、出力バッファリングを用いるほかのATMスイッチとは、すべての入力ポートからの入力セルが直接共通のバッファメモリに蓄積されるという点で、区別される。出力は、FIFO方式や単純な優先度方式に基いて選択される。このようなスイッチは、単純な優先度で満足するような、長距離テレコミュニケーションに対しては十分に機能するが、このスイッチには、ファクトリーオートメーションや、パワープラントコントロール、フルモーションビデオなどリアルタイムアプリケーションをサポートする待ち行列制御の機能はない。

【0010】

【発明が解決しようとする課題】この発明は、以上のような課題を解決するためになされたもので、スタティックなスケジューリングアルゴリズムだけでなく、ダイナミックなスケジューリングアルゴリズムに則った待ち行列制御が可能なスイッチングシステムを得ることを目的

としている。

【0011】

【課題を解決するための手段】この発明に係るスイッチングシステムは、複数の入力リンクと複数の出力リンクを持ち、各入力リンクによりネットワーク内の1以上の宛先を示す宛先情報を含むヘッダフィールドとデータとを格納したセルを受信するスイッチと、上記ヘッダフィールドの宛先情報を、各ビットが各出力リンクに対応しておりそのビットが有意味状態である時そのセルが最終宛先に到着するようにそのビットに対応する出力リンクに対してそのセルを出力することを示している宛先ビットのベクトルに変換する変換手段と、各セルに対して、上記宛先ビットのベクトルを付加するベクトル付加手段と、上記ベクトルを付加されたセルを1つ以上の出力リンクに接続する接続手段とを備え、上記接続手段は、上記ベクトルを付加されたセルを到着順に格納するキューイング手段と、キューイング手段に格納されたセルの中から到着順とは異なる順に各出力リンクに対して各出力リンクに対応する宛先ビットが有意味状態のセルを選択する選択手段とを備え、上記選択手段は、セルの出力をスケジューリングするスケジューリング手段を有し、上記スケジューリング手段は、選択順を決定するアルゴリズムを記憶する手段と、スイッチの動作に先だってそのアルゴリズムを設定する手段と、スイッチの動作中にそのアルゴリズムを再設定する手段と、上記各出力リンクに対してセル送信時間内に上記アルゴリズムを実行する手段を有しており、上記スイッチは、更に、上記出力リンクに選択したセルを出力する出力手段を備えたことを特徴とする。

【0012】上記出力手段は、上記ベクトルが除去されたセルを出力することを特徴とする。

【0013】上記アルゴリズムは、ダイナミック・プライオリティ・スケジューリング・アルゴリズムと、スタティック・プライオリティ・スケジューリング・アルゴリズムと、ラウンド・ロビン・アルゴリズムと、これらの組み合わせのいずれかから選択されることを特徴とする。

【0014】上記スケジューリング手段は、各出力リンクにセルを出力する順番を決定するために設定されたスケジューリングアルゴリズムに基づいて、キューイング手段の中に格納されたセルの出力順序を示す値をタグ値として計算し、このタグ値を有するタグをセルに付加するタグ付加手段を有し、上記選択手段は、タグ値と宛先ビットの値に基づいて、キューイング手段の中に格納されたセルをサーチするサーチ手段とを備え、選択手段は、各出力リンクに対応する宛先ビットの値が有意味状態であるセルの中で最も小さいタグ値を持つセルをサーチして選択し、もし、最も小さいタグ値を持つセルが2以上存在する場合に、先に到着したセルを選択することを特徴とする。

【0015】上記キューイング手段は、先頭から末尾に

至る複数のレジスタからなるFIFO回路を備え、各レジスタは、出力を待っているセルを示すものであり、各レジスタは、上記宛先ビットのベクトルとタグ値のバイナリーコードを保持し、各レジスタは、先頭から末尾の方向でセルの到着順を示していることを特徴とする。

【0016】上記サーチ手段は、上記レジスタの数と同じ数の複数の比較回路を備え、各比較回路は、各レジスタに対応して設けられ、各比較回路は、対応したレジスタのタグ値と末尾側にある次のレジスタに対応した比較回路からの出力値を比較するようにリニアに配置され、各比較回路は、対応するレジスタのタグ値が末尾側にある次に比較回路から出力されたタグ値より小さい場合であって、かつ、選択した出力リンクに対する宛先ビットが有意味状態である場合に、出力ビットをオンにすることにより、上記比較回路は、選択した出力リンクに対して最小のタグ値を伝搬させ、待ち行列の中からセルを選択することを特徴とする。

【0017】上記サーチ手段は、上記レジスタの数の半分に当たる複数の比較回路と追加の比較回路を備え、上記比較回路を階層的に配置し、ルートと枝と葉を持つツリーを構成し、葉と枝のサブセットによりサブツリーを構成し、葉に配置された各比較回路は、2つの隣り合うレジスタのタグ値を比較して小さいタグ値を有するレジスタを識別し、上記追加の比較回路は、2つのサブツリーからのタグ値を比較し、選択した出力リンクに対して、最小のタグ値を持つレジスタが存在しているサブツリーを識別し、上記ツリーは、選択した出力リンクに対して、最小のタグ値を伝搬することを特徴とする。

【0018】また、この発明に係るスイッチングシステムは、複数の出力リンクと出力リンクを示すアドレスを含んだセルを受信する複数の入力リンクと、上記アドレスを宛先ビットに変換する変換手段と、上記入力リンクのセルを少なくとも1つの出力リンクに接続する接続手段とを備え、上記接続手段は、入力リンクから出力リンクへセルの出力をスケジューリングするスケジューリング手段を備え、スケジューリング手段は、各出力リンクにセルを出力する順番を決定するために、予め設定されたスケジューリングアルゴリズムに基づいて、どのセルを出力すべきかを示す値を生成し、生成した値をタグ値として有するタグをセルに付加するタグ手段と、共通キューと、上記共通キューに、各セルに対して宛先ビットとタグ値を挿入する手段と、上記宛先ビットとタグ値に基づいて、上記共通キューから、対応している宛先の中から最小のタグ値を持つセルをサーチするサーチ手段と、宛先に対応した出力リンクに対して最小のタグ値を持つセルを出力する出力手段とを備えたことを特徴とする。

【0019】上記サーチ手段は、最小のタグ値をシリアルに評価していくリニアサーチ手段を有していることを特徴とする。

【0020】上記サーチ手段は、木構造を用いて同時に

複数対のセルを評価していくロガリズムックサーチ手段を有していることを特徴とする。

【0021】上記サーチ手段は、リニアサーチ手段とロガリズムックサーチ手段を結合したサーチ手段を有していることを特徴とする。

【0022】上記タグ手段は、先頭から末尾までのキューを形成するように配置された一連のタグレジスタと、複数の比較器と、上記タグレジスタと比較器を接続する接続手段と、上記セルに対応するタグレジスタにより識別できるように記憶する記憶手段と、記憶されたセルのタグ値をタグレジスタに設定する設定手段とを備え、上記比較器は、タグレジスタのタグ値と1つ前の比較器からの出力値を比較して小さい方の値を決定するとともに、小さい方の値と宛先ビットを出力する出力手段と、小さい方の値と宛先ビットに基づいて、記憶されたセルを識別する識別手段とを備えたことを特徴とする。

【0023】上記比較器は、一連のビットから構成された出力を持ち、この一連のビットの中に、宛先ビットが有意状態であり、同じ宛先ビットが有意状態であるキュー中の他のセルのタグ値よりも小さいタグ値を持つことを示す特定ビットを有しており、各比較器は、各タグ値とキューの末尾方向にあるタグの最小値とを比較することにより、宛先別に各セルのプライオリティを設定することを特徴とする。

【0024】上記比較器は、2つのタグ値を比較する手段と、小さい方の値を出力する手段とを有し、1つのタグ値に対して1つの追加の入力ビットを備え、上記比較器の出力は、比較対象となるタグ値を示す追加の入力ビットにより条件付けられていることを特徴とする。

【0025】上記追加の入力ビットは、宛先ビットが有意状態になっているセルのタグ値を特定するものであることを特徴とする。

【0026】上記スイッチングシステムは、更に、最小のタグ値を持つセルが2以上存在する場合に、キューの先頭に近いセルを選択する手段を備えたことを特徴とする。

【0027】上記タグ手段は、枝を経由してルートに連結された葉を持った階層的木構造を有し、各枝は、一対のタグレジスタと隣り合うタグレジスタの間に設けられたタグレジスタ用比較器と、2つのタグレジスタ用比較器からの出力を入力する追加の比較器を有しており、上記追加の比較器は、宛先ビットが有意状態になっており、小さい方のタグ値を持つタグレジスタに対応するセルを識別する手段を有していることを特徴とする。

【0028】上記タグレジスタ用比較器は、小さい方のタグ値と2つの比較ビットを出力し、2つの比較ビットの内、一方の比較ビットは一方のセルが小さい方のタグ値を持っていることを示し、2つの比較ビットの内、他方の比較ビットは他方のセルが小さい方のタグ値を持っていることを示すものであり、上記選択手段は、小さい

方のタグ値を持っているセルを識別するために、宛先ビットと比較ビットを入力する複数のANDゲートを有し、複数のANDゲートは、ルートにおいてただ1つのANDゲートが有意な出力を有するように構成されており、このANDゲートの有意な出力により出力リンクに出力されるセルを示すことを特徴とする。

【0029】上記ANDゲートは、そのANDゲートに対応したセルが適切な宛先ビットを有しているかを判定するために用いられるものであり、上記タグレジスタ用比較器は、異なるANDゲートへ接続された出力ビットを有し、このタグレジスタ用比較器からの出力ビットは、対応するセルの宛先ビットとANDを取られ、その結果の出力が有意状態である場合は、対応するセルが宛先ビットが有意状態であり、最小のタグ値を持つセルであることを示すことを特徴とする。

【0030】上記スイッチングシステムは、更に、2以上のセルが最小のタグ値を有する場合、先着のセルを識別する手段を有することを特徴とする。

【0031】上記タグレジスタ用比較器は、先頭と末尾を持つキューを定義し、上記追加の比較器は、2入力から小さい値を出力するとともに、キューの先頭方向にある2入力の値がキューの末尾方向にある2入力の値以下である場合に、有意状態を出力することを特徴とする。

【0032】上記スイッチングシステムは、更に、上記木構造のルート方向への上位層レベルにおいて、下位層レベルのどの枝が小さい方のタグ値を有しているか決定することにより小さい方のタグ値を有する枝を識別する識別手段を有し、上記識別手段は、各タグに対応してそのタグの宛先ビットとそのタグが最小のタグ値を有することを示すビットとのANDを取るためのANDゲートを有していることを特徴とする。

【0033】

【作用】この発明のスイッチングシステムは、スケジュール手段がセルの出力をスケジュールする。スケジュール手段は、所定のアルゴリズムに基づいて、セルを順番に選択する。このように、スケジュール手段を備えることにより、セルの優先処理が行える柔軟なシステムを提供する。また、上記スケジュール手段は、いろいろな種類のアルゴリズムを採用することができる。そのアルゴリズムによって、優先度の高いセルをタイムリーに交換することができるとともに、優先度の低いセルを後から交換することができる。

【0034】出力手段は、セルを出力する場合には、スイッチの内部で用いた宛先ビットのベクトルを除去して出力する。

【0035】上記アルゴリズムとしては、ダイナミックなスケジューリングアルゴリズムやスタティックなスケジューリングアルゴリズムを用いることができるとともに、ラウンド・ロビン・アルゴリズムも用いることができる。或いは、これらのアルゴリズムの組み合わせでも

よい。

【0036】スケジュール手段は、優先順位を示すタグ値を計算し、このタグ値を有するタグをセルに付加する。サーチ手段は、このタグ値と宛先ビットを用いて出力すべきセルをサーチする。結果として、選択手段は、ある出力リンクに対する宛先ビットが有意状態（オン）になっており、かつ、それらの中で最も小さいタグ値を持つセルを選択する。最も小さいタグ値を持つセルが2以上存在する場合には、先に到着したセルを選択する。

【0037】キューイング手段は、複数のレジスタによりFIFOを構成している。この複数のレジスタからなるFIFOにより到着順を示すとともに、宛先ビットのベクトルとタグを保持するようにしている。

【0038】サーチ手段は、複数の比較回路をリニアに配置し、各レジスタに記憶されたタグ値を順に比較することにより、最小のタグ値を持つセルを検出する。

【0039】或いは、サーチ手段は、2個のレジスタに対して1つの比較回路を備え、2個のレジスタのタグ値の比較結果を更に追加の比較回路で比較することにより、ツリー形式を用いて、最小のタグ値を持つセルを選択する。このように、ツリー構造を用いた比較処理を行うことにより、キューに記憶されたエントリーの数の対数に比例する時間で最小のタグ値を持つセルを選択でき、前述したリニアにサーチする場合に比べて高速処理を行える。

【0040】この発明のスイッチングシステムは、スケジュール手段が、セルの出力順番を予め設定されたスケジューリングアルゴリズムに基づいて決定する。スケジュール手段は、出力リンクに共通に設けられたキューを備え、そのキューに対して各セルの宛先ビットと優先度を示すタグを記憶させる。サーチ手段が、この宛先ビットとタグを参照することにより、出力リンクに出力するセルの順番を決定する。タグに記載されるタグ値は、ダイナミック、或いは、スタティックなアルゴリズムのいずれのアルゴリズムに基づくもので計算されたもので構わない。従って、このスイッチングシステムは、ユーザの要求に応じて柔軟なスケジューリング方式を採用することができる。

【0041】上記サーチ手段は、リニアサーチ手段により、最小のタグ値を順番に捜していく。

【0042】或いは、サーチ手段は、ログリズミックサーチ手段を有し、木構造を用いて最小のタグ値を高速に検索する。

【0043】或いは、サーチ手段は、リニアサーチ手段とログリズミックサーチ手段を結合したサーチ手段により、最小のタグ値を検索する。

【0044】上記タグ手段は、タグレジスタと比較器を備え、タグレジスタにタグ値を設定し、比較器がこのタグ値を順番に比較していくことにより、最小のタグ値を検出する。

【0045】上記比較器は、特定の出力リンクに対する宛先ビットが有意状態（オン）であり、最小のタグ値を持つことを示す特定ビットを出力し、各比較器は、キューの末尾方向にある比較器からの出力と各比較器に対応したタグ値とを比較することにより、優先度を決定する。

【0046】比較器は、2つのタグ値を比較して小さい方の値を出力する。また、比較器は、追加の入力ビットを備え、追加の入力ビットにより入力したタグを比較対象とするかどうかを決定する。

【0047】上記追加の入力ビットは、宛先ビットが有意状態（オン）になっているセルを示すことにより、対応するタグ値が比較対象になるかどうかを特定するものである。

【0048】もし、最小のタグ値を持つセルが2以上存在する場合には、先に到着したセルを選択する。

【0049】前述したタグ手段の別な構成として、タグを階層的構造を持った比較器により比較するようにしても構わない。一对のタグレジスタに対して、比較器を1個設け、この比較器からの出力を更に追加の比較器で比較することにより、最小のタグ値を持つタグレジスタを検出する。このようにして、検索を高速に行う。

【0050】比較器は、比較した2つのタグ値の内、いずれのほうか小さい値を持っているかを示す比較ビットを出力し、選択手段は、複数のANDゲートを有し、ANDゲートが比較ビットを入力するとともに、宛先ビットを入力することにより最終的に優先度の高いセルを1つだけ特定する。

【0051】ANDゲートは、宛先ビットが有意状態（オン）になっているかどうかを判定するとともに、比較器からの比較ビットとのANDを取ることに、最終的に宛先ビットが有意状態（オン）であり、最終のタグ値を持つセルを選択する。

【0052】もし、2以上のセルが最終のタグ値を有する場合には、先に到着したセルを選択して出力する。

【0053】タグ手段が前述したように、階層的木構造を用いて最小のタグ値を検索する場合には、タグレジスタ用比較器により先頭から末尾方向が存在するキューを定義し、追加の比較器により先頭方向にあるタグ値が末尾方向にあるタグ値よりも小さい場合に、その先頭方向にあるタグ値を持ったセルを選択する。

【0054】本発明のスイッチングシステムは、更に、識別手段により階層的木構造を持つ場合に、どの枝が小さい方のタグを有しているかを識別し、その識別手段は、ANDゲートに基づいて、タグの宛先ビットとそのタグ値が最小のタグ値であることを示すビットとのANDを取ることに、最終的に優先度の高いセルを選択する。このようなANDゲートを備えることにより、階層的木構造を取った場合に、葉からルートに対して宛先ビットを伝搬するため、各タグ毎に必要とされていたA

NDゲートを省略することができる。結果として、ANDゲートの数を減少させることができ、回路構成を簡略することができる。

【0055】

【実施例】

実施例1. 各種のスケジューリングアルゴリズムを適合させるために、この発明に係るデジタル通信ネットワークスイッチは、改良されたスケジューリングシステムを備えている。その改良されたスケジューリングアルゴリズムによれば、各セルは、スイッチに到着する都度、10 タグ付けされ、その後、共通キューに記憶される。タグは、バイナリーの数値から成る。タグの数値は、トラヒックのクラスに対応するスケジューリングアルゴリズムや、セルが伝送されるバーチャルチャネル(virtual channel)の属性や、セル自身の特性を考慮して計算される。続いて、キューが宛先とタグにより出力リンク又は出力ポート毎に並列に検索される。複数のセルの宛先が同一である時、即ち、出力リンク又は出力ポートが同一である時は、その中で最小のタグ値を持つセルが選択される。これにより、スイッチは、常に、20 宛先毎に、スケジューリングアルゴリズムに従って、最初のセルを選択する。

【0056】より一般的に言えば、この発明のスイッチングシステムは、アルゴリズムを記憶する手段を備えている。記憶されたアルゴリズムに従って、セルがスケジュールされる順番が決定される。また、アルゴリズムは、スイッチの動作に先立って、また、スイッチの動作中に変更可能であるので、ネットワークの幅広いトラヒック要求に対応することが可能となる。

【0057】この発明のスイッチは、スタティックアルゴリズムとダイナミックアルゴリズムの両方をサポート可能である。これらのアルゴリズムには、以下のものを含む。ルー、ジャングによるレートモニタリングアルゴリズム、ジャングによる単純優先度アルゴリズム、クマーによるウェイトドフェアキューイングアルゴリズム、ジャクソン、ルー、ホーン、ジェングによる最早デッドラインファーストアルゴリズム、ジャングによるバーチャルクロックアルゴリズム、カルマネクによるラウンド・ロビン及び階層的ラウンド・ロビン。以上のアルゴリズムは、本明細書で既に述べたものである。また、40 カングのフローコントロールドバーチャルチャネルアルゴリズムのような複雑な制御アルゴリズムもサポート可能である。

【0058】この発明は、以下の3点により成り立っている。第1に、タグ値がいずれかの演算可能なファンクションに従って計算される点、第2に、タグ値がセルが伝送される順番を示している点、第3に、上に挙げたすべてのアルゴリズムに共通する特徴として、セルの順番が数値のタグで表現されるという点である。検索を成立させる唯一の決定要因は、タグの持つ数値であり、タグ

の持つ意味は問われない。さらに、各セルのタグ値は、伝送中に計算されるので、スタティック・スケジューリング・アルゴリズムだけでなく、ダイナミック・スケジューリング・アルゴリズムもサポートできる。また、タグ毎に、十分な数のビットが用意されているので、スイッチバッファ中の全セルの順番を並び替えて、各ビットの値を書き変えることも可能である。

【0059】この発明のポイントは、タグベースの検索を行う点にある。各セル(ATMセル)は、ネットワークスイッチに到着すると、バイナリーの数値を持つタグを付加される。典型的な実施例においては、8ビット、16ビット、あるいはそれ以上の複数のビットを持つタグが用いられる。8ビットの場合には、256(2の8乗)通りのタグ値が設定できる。16ビットの場合には、65536(2の16乗)通り、nビットの場合には、2のn乗のタグ値が設定できる。各セルのタグは、そのセルが伝送されるバーチャルチャネル、そのバーチャルチャネルのトラヒックのクラスのスケジューリングアルゴリズム、そのセル自身の特性等に関連する情報を持っている。そして、セル、タグ及び宛先情報が、伝送待ちのセルキューの末尾に入力される。この実施例においては、キューは、VLSI回路で構成されている。

【0060】やがて、スイッチが、特定の宛先にセルを伝送しようとする時、キューに記憶されているすべてのセルを対象に出力リンク毎、或いは、出力ポート毎に並列に検索が行われ、その宛先に一致する最小のタグ値を持つセルが選択される。選択されたセルは、キューから消去され、伝送される。キューの中に同一のタグ値を持つ複数のセルが存在する場合には、もっとも早くキューに到着したセルが選択される。キュー検索部は、最小のタグだけでなくセルが伝送される出力ポートも識別する。

【0061】このアーキテクチャにより、スタティック・スケジューリング・アルゴリズム、ダイナミック・スケジューリング・アルゴリズムを問わず、前述した全てのアルゴリズムを含む幅広いクラスのスケジューリングアルゴリズムを実現することが可能となる。また、タグ値に範囲設定をすれば、1つのスイッチで同時に種々のトラヒックのクラスに対応する種々のスケジューリングアルゴリズムを実現できる。具体的には、タグ値を計算する際に、優先度の高いクラスのタグを数値の小さい範囲に割り当て、優先度の低いクラスのタグを数値の大きい範囲に割り当てる。そのような割り当てを行えば、数値を頼りに検索するシステムでは、優先度の高いクラスのセルが1つでもあれば、必ずそのセルが選択される。なぜならば、優先度の高いクラスのセルのタグは、優先度の低いクラスのセルのタグよりも、より小さい値を持つからである。

【0062】図1を用いて説明する。ATMスイッチ10は、複数の入力リンク14を備えている。複数の入力

17

リンク14は、複数の入力処理部16に接続されている。複数の入力処理部16からの出力は、それぞれキュー検索部18とメモリ（セルバッファメモリ）20に入力される。キュー検索部18は、制御部22により制御される。制御部22は、また、入力処理部16及び出力処理部24も制御する。入力処理部16及び出力処理部24は、それぞれマイクロプロセッサとメモリを備えている。また、入力処理部16と出力処理部24は、制御部22によりアルゴリズムを設定、或いは、再設定する手段を備えており、スイッチに到着するセルやスイッチから送り出されるセルのネットワークへの要求に従って、適切なアルゴリズムが用いられるように制御される。

【0063】キュー検索部18の出力は、それぞれ出力処理部24に入力される。それに伴い、セルバッファメモリ20の出力も出力処理部24に入力される。26に示すように、入力処理部16の出力は、宛先情報とタグとバッファアドレスを含んでいる。これらは、キュー検索部18に接続されている。また、28に示すように、セルヘッダーとセル本体は、セルバッファメモリ20に記憶される。

【0064】次に、動作について説明する。セルは、入力リンク14を介して入力処理部16に到着する。入力処理部16は、到着したセルをセル毎に処理する。通常、どのATMスイッチでも行われている内部処理に加えて、入力処理部16はセル毎にタグ値を計算する。タグ値の計算は、そのセルのバーチャルチャネルに対応するスケジューリングアルゴリズムに従って行われる。タグ値は、また、バーチャルチャネルの状態やセルの到着時刻にも左右される。バーチャルチャネルとは、他の全てのデータストリームと区別される特定のデータストリームのエンドトゥーエンドのコネクションを指す。ATMネットワークにおいて、それは、セルヘッダー中のビットにより識別される。入力処理部16は、また、バーチャルチャネルの状態情報やスイッチによりメンテナンスされる他の情報も計算し、更新する。

【0065】上記計算後、タグと宛先情報とセルのバッファアドレスとが出力され、キュー検索部18のキューに記憶される。また、セルヘッダーとセルデータは、セルバッファメモリ20に出力される。その時点で、セルのアドレス、宛先及びタグはキューイングされ、セルヘッダーとセルデータは、セルバッファメモリに記憶される。こうして、入力処理部は、次のセルを受け付け、処理することが可能な状態になる。以上のことから、入力処理は、関連する入力リンク14の1セルサイクルに割り当てられた時間内に完了しなければならない。1セルサイクルとは、入力リンクの帯域幅で1つのATMセルを受信するのに必要な時間のことである。

【0066】キュー検索部18は、キューを順に検索し、個々の出力処理部24について対応する宛先ビット

18

を持つセルを選択する。出力処理部24は、セルが伝送されるのに必要な計算があればその計算を行う。その後、ネットワークの他のノードやセルの最終宛先に向けて出力リンク30にセルを伝送する。キュー検索部18によって行われる検索は、全ての出力処理部に関連するスイッチ内の入力リンク14及び出力リンク30の内、最も遅いリンクの1セルサイクルに割り当てられた時間内に完了しなければならない。

【0067】この発明のキュー検索システムの説明に先だって、図2を用いて既に知られている検索システムについて説明する。これは、図1に示すキュー検索部18に相当する。図2に示す検索システムは、縦列接続することにより上位ステージ又は／及び下位ステージを備えることができ、拡張することも可能である。キュー検索部18は、VLSIを用いたFIFO回路31（self-timed FIFO）からなっている。キューエントリーは、FIFO回路31の末尾から挿入される。FIFOに未使用のエントリーがあれば、自動的にキューの先頭方向に前詰めにシフトされる。各キューエントリー39は、レジスタに記憶され、宛先フィールド41とアドレスフィールド35と優先度ビット37を持っている。宛先フィールド41は、独立した一連のバイナリービットを持っており、各ビットは1か0かいずれかの値を持つ。この実施例では、1の値を有意味状態（オン）としている。この値は、データがどの出力リンクに伝送されるかによって決定される。宛先フィールドを構成している各ビットを、宛先ビットと呼ぶ。スイッチの出力リンク30と同数の宛先ビットがあり、各宛先ビットは、特定の出力リンクへのコネクションを示している。キューエントリーの宛先フィールドの対応するビットが1の時に、キューエントリーと出力リンクを結ぶコネクションが成立する。例えば、宛先フィールドの第3ビットが1の時には、対応するデータは、第3の出力リンクに伝送されることを示している。

【0068】上で述べたように、宛先フィールドの各ビットは、それぞれ出力リンクに対応しているので、キュー全体を縦の列に着目して先頭から末尾まで眺めてみると、キューに記憶している全てのセルと1つの出力リンクとの対応関係が分かる。即ち、各カラム（列）のビットが1か0かを調べると、セルが特定の出力リンクに伝送されるかどうかを確認することができる。マルチキャストセルの場合、即ち、1つのセルが複数の出力リンクに関連付けられている場合には、そのセルの宛先フィールドには、1の値を持ったビットが複数存在する。図2において、例えば、宛先フィールドの右端のカラムを見てみると、1の値を持ったセルが2つあることが分かる。この従来の検索システムにおいては、どのスケジューリングアルゴリズムを用いても、同一の出力リンクに向かう2つのセルがある場合に、その伝送の順番を決定できるようにするために、優先度ビット37を設けてい

る。

【0069】以上のように、キューエントリー39の宛先フィールド41の中に1の値を持つビットがあれば、そのキューエントリーのアドレスフィールド35に示されたアドレスに対応して、セルバッファメモリ20に記憶されているセルは、必ずそのビットに対応する出力リンクに伝送される。このことは、宛先ビットによってセルが1つの出力リンクに伝送されるか、或いは、複数の出力リンクに伝送されるかを識別できることを示している。図2において、43に示すキューエントリーのように、宛先フィールド41の全てのビットが0である場合には、このエントリーは未使用であるとみなされる。また、図2には、検索回路50とパス読み出し回路45が示されている。

【0070】次に、動作について説明する。制御部22からキュー検索部18に対して、「ある1つの出力リンク30に伝送するセルを検索せよ」という命令が出されると、宛先フィールド41の中でその出力リンクに対応するビットが特定される。次に、検索回路50が特定されたビット位置に1がたっていて、かつ、キューの先頭に最も近いキューエントリー39を選択する。次に、制御部22は、選択されたキューエントリーのセルのアドレスフィールド35をパス読み出し回路45に入力する。パス読み出し回路45において、選択されたセルのアドレスフィールド35は、キュー検索部18の出力49となる。続いて、選択されたキューエントリーの対応する宛先ビットは、リセットされ、値は0となる。最後に、出力処理部24は、アドレスを出力49により受け取る。そのアドレスは、セルバッファメモリ20から指定された出力リンク30へ伝送するセルを選択するために用いられる。

【0071】図3は、説明を簡単にするために、図2に示すキューエントリー39が4つある場合を示している。また、図2において、宛先フィールド41で示していたものの内、宛先ビットを42として示している。また、遅延優先度フィルタ48は、優先度ビット37を入力する。キュー検索部18は、出力リンク30に対応するセルを選択するコマンドを受けると、宛先ビットの中からその出力リンクに対応するカラム46を選択する。このカラムは、遅延優先度フィルタ48に入力される。遅延優先度フィルタ48には、既にキューの各エントリーに対応する2つの優先度を示す優先度ビットが入力されている。出力リンクが選択された時、宛先ビットの内、その出力リンクに対応するカラムの全てのセルに対応するビットが、遅延優先度フィルタ48に入力される。遅延優先度フィルタ48は、優先度ビットと宛先ビットを結合し、結合結果として0か1の値を持つ各セルに対応するビットを出力する。この時、出力された値が0であれば、そのセルは選択されないことを示している。また、ビットの値が1である時には、そのセルが選択され

る資格があるということを示している。このように、遅延優先度フィルタ48は、2種類の優先度を用いて単純な優先度付けの機能を果たしている。この従来技術において、遅延優先度フィルタは、同じ宛先を持つ複数のセルの中から、いくつかの任意のセルを選択するために用いられている。より明確に言えば、遅延優先度フィルタは、高い優先度と低い優先度をつけることによって、高い優先度のセルが低い優先度のセルよりも、先に伝送されることを可能にしている。

10 【0072】遅延優先度フィルタ48から出力された一連のビットは、比較回路50aに入力される。比較回路50aは、各エントリーのビットをその両隣のビットと比較する。これは、図3に示されている論理和ゲート52及び排他的論理和ゲート54により行われる。その結果、出力ビットの内、1つのビットだけが1の値を持つことになる。それは、その1の値を持つビットに対応するエントリーが選ばれて伝送されることを示している。選択されたエントリーは、宛先ビットが1であるエントリーの内、最もキューの先頭に近いエントリーである。

20 【0073】図3に示すアクセラレータ60は、検索をスピードアップするために用いられる。その結果、1セルサイクル内で全ての出力リンクの検索を順に行うことが可能となる。アクセラレータ60は、通常一般的に用いられている先読み回路であり、縦列接続される場合の下位ステージからの入力62と上位ステージへの出力64を備えている。下位ステージからの入力62が0の時は、下位ステージにおいて、指定された出力リンクに対してエントリーが未選択であることを示している。下位のステージからの入力62が1の時は、下位のステージにおいて、指定された出力リンクに対してエントリーが選択されたことを示している。従って、下位ステージからの入力が1の時は、比較回路からの出力ビット56は、全て0になるように設計されている。上位ステージへの出力64は、下位ステージ又は自ステージにおいて、エントリーが未選択の場合、0が出力される。下位ステージ又は自ステージにおいて、エントリーが既に選択された場合は、1が出力される。

40 【0074】図4は、図3における従来の遅延優先度フィルタ48を、この発明のリニアサーチ回路70で置き換えたものである。この実施例においては、各キューエントリー39は、タグレジスタ72により拡張されている。タグレジスタ72は、数値に対応する値を持っており、その数値の大きさは、優先度を設定するために用いられる。この発明においては、最も小さいタグ値を持つセルが最も高い優先度を与えられる。前述したように、タグの持つ値によって優先度が設定される。この実施例のリニアサーチ回路70は、様々な選択機能を提供することを目的として備えられている。これは、前述した図3に示す遅延優先度フィルタ48が単純化機能を目的としていた点と大きく異なっている。多様な選択機能を提

21

供することにより、複数のスケジューリングアルゴリズムの幅広い要求に柔軟に答えることができる。図4においても、説明を簡単にするために、数百のキューエントリーの中から、4つのエントリーを取り分けている。4つのエントリーに対して、それぞれタグレジスタ72と比較器74が対応している。各タグレジスタ72は、図1に示す入力処理部16で計算された、セルエントリーに対応するタグの値を示すバイナリーの数値を持っている。この数値が比較器74において、キューの中の直前のタグの比較器の出力と比較される。比較結果により、2つの数値の内、どちらが小さいかが決定される。この順序付け処理の結果は、キューの中を順に次の比較器に引き継がれる。

【0075】また、この順序付け処理の結果は、一連のビットとして出力ライン76に出力される。これらの一連のビットの内、いずれかのビットに1がたっていれば、それは第1に、対応するキューエントリーの宛先ビットが1であることを示している。第2に、対応するキューエントリーのタグ値がキューの中で選択された宛先ビットが1である全てのキューエントリー中で、そのキューエントリーのタグ値が最も小さいということも示している。それ故、以下のことが分かる。1つは、リニアサーチ回路70及び比較回路50a及びアクセラレータ60からなるリニアサーチ回路の出力ビット56は、選択された宛先ビットが1にセットされている全てのキューエントリーの中から、1つだけを選択するということ、もう1つは、選択されたキューエントリーは、最小のタグ値を持つということである。また、更に、もしも2つ以上のキューエントリーが同一の最小のタグ値を持ち、同一の宛先ビットを持っている場合に、キューの先頭に最も近いエントリーが選択されるということになる。

【0076】その結果、リニアサーチ回路70により、スケジューリングアルゴリズムが適当であろうと思って採用した方式に従って計算されたタグ値に基づいて、選択プロセスを独自に構築する機会が提供される。即ち、リニアサーチ回路70は、タグ値が小さな値を持つセルを先に出力するように設計されているのであるから、スケジューリングアルゴリズムが、先に出力したいセルのタグ値を小さな値にすることさえ守れば、どのような方式でタグ値を決定してもよいことになる。このようにして、スケジューリングアルゴリズムの独自構築が可能になる。

【0077】特に、この発明においては、タグベースの検索を実現するために、従来の遅延優先度フィルタ48を使用していない。遅延優先度フィルタ48の代わりに、キューの各エントリーに対応するタグレジスタを持つ。また、各タグレジスタをキューの末尾に向かって他のタグと比較し、最小値を持つタグを選択する比較器74を持つ。図5は、図4に示したりニアサーチ回路70

22

の詳細を示す図である。図5に示すように、以下の例において、 k を各タグレジスタのビット数とする。また、 N をキューエントリーの総数とする。キューエントリー i に対応するタグをタグ i とし、キューエントリー i に対応する比較器を比較器 i とする。タグ i のタグ値をタグ値 T^i とする。キューエントリー i に対応する宛先ビットを宛先ビット D^i とする。比較器 $i+1$ と比較器 i の間を通過している一組のワイヤの出力値を出力値 M^{i+1} とする。比較器 i の出力を、出力値 M^i 、1ビット C^i とする。この1ビット C^i は、宛先ビット D^i がセットされているかどうかとともに、もし、セットされている時は、更に、タグ値 T^i が出力値 M^{i+1} 以下であるかどうかを示す。

【0078】キュー最後尾エントリーからの出力値 M^N は、宛先ビット D^N に1がたっている場合には、タグ値 T^N の値にセットされる。また、宛先ビット D^N が0の時は、 $2^k - 1$ （全ビット=1）という値になる。比較器 i の一組のワイヤからの出力値 M^i は、キューエントリー i からキューの最終（キューエントリー N ）までのエントリーの内、最小のタグ値を示す。また、ビット C^i は、宛先ビット D^i がセットされ、かつ、そのタグ値 T^i がキューエントリー i から N までのキュー内にあるタグ値の中で最小値の場合、1に設定される。出力されたビット C^i は、図2及び図3に示す従来例の回路の遅延優先度フィルタの出力を置き換える。検索回路50は、単に、ビット C^i が1にセットされている最初のエントリーを選択する。即ち、最小のタグ値を持つエントリーを検索する。

【0079】以下に、比較器の動作について、詳細に説明する。タグレジスタの各ビットを比較する比較論理回路を図6に示す。図6は、図5に示す比較器 i の j 番目の比較論理回路を示している。図6の $T^i[j]$ とラベル付けされたボックス130は、 i 番目のキューエントリーのタグレジスタの j 桁目（ $k-1 \geq j \geq 0$ ）のビットを示している。 $M^{i+1}[j]$ とラベル付けされたワイヤ132は、出力値 M^{i+1} の j 番目のビットであり、キュー最後尾方向にある $i+1$ 番目のキューエントリーの比較器 $i+1$ からくるものである。また、 $M^i[j]$ とラベル付けされたワイヤ132は、出力値 M^i の j 番目のビットであり、このエントリーからキューの先頭に向かって出力されるものである。 $C^i[j+1]$ とラベル付けされたワイヤ136と、 $E^i[j+1]$ とラベル付けされたワイヤ138は、 i 番目のエントリーの比較器 i の（ $j+1$ ）番目の比較論理回路からくるキャリービットである。また、 $C^i[j]$ とラベル付けされたワイヤ140と $E^i[j]$ とラベル付けされたワイヤは、 i 番目のエントリーの比較器 i の j 番目の比較論理回路から次の下位ビットの比較論理回路へ、即ち、（ $j-1$ ）番目の比較論理回路へ出力されるキャリービットである。

23

【0080】2つのキャリービットとワイヤからの出力値の意味は、以下のように定義される。

$$\begin{aligned} C^i[j] &= 1 \\ &\text{if and only if } D^i = 1 \\ &\text{and } T^i[(k-1) \dots j] \\ &\leq M^{i+1}[(k-1) \dots j] \end{aligned} \quad (1)$$

$$\begin{aligned} E^i[j] &= 1 \\ &\text{if and only if } D^i = 1 \\ &\text{and } T^i[(k-1) \dots j] \\ &= M^{i+1}[(k-1) \dots j] \end{aligned} \quad (2) \quad 10$$

$$\begin{aligned} M^i[(k-1) \dots j] \\ &= \min \{ 2^k - 1, \\ &\quad T^m[(k-1) \dots j] \\ &\quad \text{for } m=1, \dots, N, \\ &\quad \text{such that } D^m = 1 \} \end{aligned} \quad (3)$$

ここで、N=キューのエントリー総数、 $N \geq 1 \geq 1$ 、 $k-1 \geq j \geq 0$ である。また、 $T^i[(k-1) \dots j]$ は、kビットで示されたタグ値 T^i の上位 $k-j$ ビットの値を意味する。

【0081】式(1)により、 $C^i[j]=1$ の時は、宛先ビット D^i が1であり、かつ、エントリーiのタグ値 T^i の上位 $k-1$ ビットからjビットまでの値がエントリー $i+1$ からエントリーNまでの最小のタグ値 M^{i+1} の上位 $k-1$ ビットからjビットまでの値以下であることを示している。 $C^i[j]=0$ の時は、そのエントリーiは、検索すべきエントリーでないことを示す。 $C^i[0]$ は、 C^i として比較回路に出力される。また、式(2)により、 $E^i[j]=1$ の時は、 $C^i[j]=1$ の時で、かつ、タグ値 T^i と出力値 M^{i+1} の上位 $k-1$ ビットからjビットまでの値が等しいことを示している。従って、 $E^i[j]=1$ の時は、必ず $C^i[j]=1$ である。式(3)は、 $M^i[(k-1) \dots j]$ がエントリーiからエントリーNまでの最小のタグ値の上位 $k-1$ ビットからjビットまでであることを示している。

【0082】図6に示す比較論理回路の機能は、図7に示す論理値表で定義される。図7において、5行目の $C^i[j+1]$ と $E^i[j+1]$ の値は、ありえない組み合わせである。図7を用いて、比較器が最小値を導き出すことを以下に証明する。まず、最初に、以下のよう

$$\begin{aligned} C^i[k] &= E^i[k] = D^i \\ M^i[j] &= T^i[j] \quad \text{or} \quad (-D^i), \\ &\quad \text{for } j=0, \dots, k-1 \end{aligned}$$

【0083】式(1)及び(2)は、全てのi及びj=kについて成り立つ。同様に、式(3)もi=N及び全てのj=0, ..., k-1について成り立つ。次に、帰納法で話を進めよう。 $0 \leq j < k$ 、 $1 \leq i < N$ の時のビット $T^i[j]$ を考えてみる。上記全ての値及び図6

24

のi, j以外の値について、仮定を満足したとする。すると、この仮定により、

$$\begin{aligned} C^i[j+1] &= 1 \\ &\text{if and only if } D^i = 1 \\ &\text{and } T^i[(k-1) \dots (j+1)] \\ &\leq M^{i+1}[(k-1) \dots (j+1)] \end{aligned} \quad (4)$$

$$\begin{aligned} E^i[j+1] &= 1 \\ &\text{if and only if } D^i = 1 \\ &\text{and } T^i[(k-1) \dots (j+1)] \\ &= M^{i+1}[(k-1) \dots (j+1)] \end{aligned} \quad (5)$$

$$\begin{aligned} M^i[(k-1) \dots (j+1)] \\ &= \min \{ 2^k - 1, \\ &\quad T^m[(k-1) \dots (j+1)] \\ &\quad \text{for } m=1, \dots, N, \\ &\quad \text{such that } D^m = 1 \} \end{aligned} \quad (6)$$

となる。また、図7に示す比較器のロジックテーブルの1行目~4行目の場合には、 $C^i[j+1]=0$ であるから、 $D^i=0$ であるか、または、

$$\begin{aligned} T^i[(k-1) \dots (j+1)] \\ > M^{i+1}[(k-1) \dots (j+1)] \end{aligned} \quad 20$$

であり、1行目~4行目の5列目に示すように、

$$M^i[j] = M^{i+1}[j]$$

となる。その結果、

$$\begin{aligned} T^i[(k-1) \dots j] \\ > M^{i+1}[(k-1) \dots j] \end{aligned}$$

以上のように、式(1)、(2)及び(3)は、成立する。5行目は、ありえない $C^i[j+1]$ と $E^i[j+1]$ の組み合わせを示している。6行目~9行目の場合には、 $C^i[j+1]=1$ であり、 $E^i[j+1]=0$ なので、 $T^i[(k-1) \dots (j+1)]$ が最小値となり、

$$\begin{aligned} T^i[(k-1) \dots (j+1)] \\ < M^{i+1}[(k-1) \dots (j+1)] \end{aligned} \quad (7)$$

及び、6行目~9行目の5列目に示すように、

$$M^i[j] = T^i[j]$$

となる。その結果、以下の式が成立する。

$$\begin{aligned} T^i[(k-1) \dots j] \\ < M^{i+1}[(k-1) \dots j] \end{aligned}$$

また、6行目~9行目の6列目と7列目に示すように、 $C^i[j]$ と $E^i[j]$ の値も対応してそれぞれ1及び0にセットされる。更に、式(7)と帰納法の仮定により、以下の式が成立する。

$$\begin{aligned} M^i[(k-1) \dots (j+1)] \\ = T^i[(k-1) \dots (j+1)] \end{aligned}$$

その結果、 $m=1, \dots, N$ 、 $D^m=1$ に対して、以下の式が真となる。

25

$$\begin{aligned}
& M^i [(k-1) \dots j] \\
& = T^i [(k-1) \dots j] \\
& = \min \{2^k - 1, T^i [(k-1) \\
& \quad \dots j], M^{i+1} [(k-1) \\
& \quad \dots j]\} \\
& = \min \{2^k - 1, T^i [(k-1) \\
& \quad \dots j]\} \quad (8)
\end{aligned}$$

10 行目～13行目の場合には、 $C^i [j+1] = 1$ 、かつ、 $E^i [j+1] = 1$ なので、 $T^i [(k-1) \dots (j+1)]$ は、今までの最小値 $M^{i+1} [(k-1) \dots (j+1)]$ と同じ値を示している。従って、

$$\begin{aligned}
& T^i [(k-1) \dots (j+1)] \\
& = M^{i+1} [(k-1) \dots (j+1)]
\end{aligned}$$

となる。それ故、10行目と13行目に示すように、 $T^i [j] = M^{i+1} [j]$ の時、

$$\begin{aligned}
& T^i [(k-1) \dots j] \\
& = M^{i+1} [(k-1) \dots j]
\end{aligned}$$

となり、11行目に示すように、 $T^i [j] < M^{i+1} [j]$ の時、

$$\begin{aligned}
& T^i [(k-1) \dots j] \\
& < M^{i+1} [(k-1) \dots j]
\end{aligned}$$

となり、12行目に示すように、 $T^i [j] > M^{i+1} [j]$ の時、

$$\begin{aligned}
& T^i [(k-1) \dots j] \\
& = M^{i+1} [(k-1) \dots j]
\end{aligned}$$

となる。これにより、式(1)、(2)及び(3)が導かれる。以上のように、式(1)、(2)及び(3)は、キューに記憶されている全てのセルのタグに適合する。

【0084】その結果、エントリー1のタグ値 T^i が最小の時、

$$\begin{aligned}
& C^i [0] = C^i \\
& T^i [(k-1) \dots 0] \\
& = \min \{2^k - 1, T^i [(k-1) \dots 0] \\
& \quad \text{for } m=1, \dots, N, \text{ such} \\
& \quad \text{that } D^m = 1\}
\end{aligned}$$

となる。以上のように、キューに対して高速なリニアサーチが行われ、最小値を持つタグが選択される。

【0085】この比較回路には、例えば、ワイヤ136と138のように、2本のキャリーラインが必要である。もし、キャリーラインが1本しかない場合には、図7に示した7行目と8行目、11行目と12行目のように、 $T^i [j]$ と $M^{i+1} [j]$ が異なる場合と、6行目と9行目、10行目と13行目のように、 $T^i [j]$ と $M^{i+1} [j]$ が同じ場合との差異を検出することができない。

【0086】この実施例のように、各キューエントリーの間に使用される単純比較器をVLSIで構成しリニアサーチを行うには、1ビットにつき、およそ25個の

26

トランジスタが必要である。そのため、タグが16ビットで構成されている場合には、約400個のトランジスタが必要となる。 $M^i [0]$ の値を出力するためには、宛先ビットからの情報が必要である。この情報は、N個全てのキューエントリーのkビットのタグの各ビットを逐次、走査比較することによって得られる。この走査比較は、配列の左上から右下に向けて順次実行される。1つのキューエントリーに対する比較結果が出るまで待つことなく、次のキューエントリーの第1ビット目の走査比較を先に開始しても構わない。このように、各比較器の走査比較を可能な限り早めに行うことにより、回路内の最長パスの長さは、 $N+k$ ビットとなる。

【0087】キュー検索部を0.8μmのCMOSで構成した場合、単純比較器は、およそ16ビットのタグ比較を5.5ナノセカンドで行う。このことから、キューに16個 ($N=16$) のエントリーがあり、そのj番目のビットを走査比較するには、ほぼ同じ時間がかかることが想定できる。例えば、16ビット ($k=16$) のタグを有する128エントリー ($N=128$) を持つキューから最小値を持つエントリーを選択するには、以下の式から概略49.5ナノセカンドかかることが計算できる。

$$\begin{aligned}
& N \text{ビットの走査比較時間} + k \text{ビットの走査比較時間} = \\
& (128/16) \times 5.5 + 5.5 = 49.5
\end{aligned}$$

【0088】上記49.5ナノセカンドは非常に速いが、この値は、キューエントリー数の増加により大きくなってしまふ。例えば、キューエントリーが256個であれば、128個の場合のほぼ倍の時間がかかる。

【0089】以上のような、キューをリニアサーチする方法には、時間がかかるという課題が残っている。この課題を解決するために、図8に示すような二分木の形式を用いてもよい。この改良例においては、一組(2つ)のキューエントリーのタグの間に1つの比較器を備えている。この比較器からの出力は、次の上位レベルの比較器の片側に入力される。あるレベルの2台の比較器に対応して、必ずその上位レベルに1台の比較器が配置されている。そのため、キューエントリーされているデータ数 $N=2^k$ の時、比較器の数はリニアサーチの場合と同じく、 $2^k - 1$ となる。図8においては、比較木(比較に用いる木構造)は、3レベルの深さとなっている。

【0090】図8は、この実施例のキューエントリーが8個ある場合の比較器の配置を示す図である。モジュール80は、二分木形式で2つの比較器から構成されるモジュールである。1つのキューエントリーからのタグレジスタ82と、隣接するキューエントリーからのタグレジスタ84が比較器86に接続される。比較器86において、入力された2つのタグの値が比較され、より小さい方のタグが選択される。選択されたタグに対応するキューエントリーの宛先ビットが1にセットされる。もしも、入力された2つのタグの値が同一である場合には、

比較器86はキューの先頭により近い方のタグを選択する。このような構成を取れば、キューに記憶されている全てのエンタリーから最小の値を持つタグを選択するのに必要な時間は、キュー内に記憶されているエンタリーの数のログリズムに比例する。例えば、キューのエンタリーが128個であった場合には、最小値を持つタグを選択するのに要する時間は、図4に示すリニアサーチ回路を用いた場合に要する時間の1/2以下となる。

例1. エンタリーの数=8の時、 $\log_2 8 = 3$

例2. エンタリーの数=128の時、 $\log_2 128 = 7$

【0091】比較器86の出力及び隣接する比較器88の出力は、上位レベルの比較器90に入力される。比較器90の出力は、同様に、比較器92に入力される。以上のように、図8においては、比較器の配置が木構造をなしている。

【0092】各比較器からは、2つの出力ビットC及びC' (図中、バーCで示す) 出力される。例えば、比較器86から出力された出力ビット94は、論理積ゲート98に入力され、もう1つの出力ビット96は、もう1つの論理積ゲート100に入力され、最小値を持つタグの選択に使用される。比較器からの出力ビットは、論理積ゲート98と100で、それぞれ選択されたキューエンタリーの対応する宛先ビットとAND演算される。その結果、論理積ゲートからの出力ビットが1である場合には、以下の3点を示している。第1に、対応するキューエンタリーの宛先ビットに1の値が入っているということ。第2に、その比較器に入力された複数のタグからなるサブツリーの中で、最小のタグ値であるということ。第3に、同一の最小タグ値を持つエンタリーが2つ以上あった場合に、そのエンタリーがそれらの中で最も早くキューイングされていること、即ち、キューの先頭に近いことを示している。このように、比較器と対応する論理積ゲートを木構造で配置したことにより、該当する宛先ビットに1がたっており、かつ、キューの中で最小のタグ値を持つ最もキューの先頭に近いキューエンタリーを選択することができる。

【0093】比較器の動作についてより詳しく説明する。木構造の葉にあたる比較器は、kビットからなる2つのタグレジスタの値を入力し、一組のkビットからなる出力を生成する。出力されたkビットには、2つの入力の内、より小さい方の値が示されている。また、木構造の葉にあたる比較器の2つの出力ビットCとC'の内、出力ビットCに着目してみると、比較器への2つの入力の内、キューの先頭側からの入力値が、キューの最後尾側からの入力値よりも小さい時、或いは、同一の値である時のみ、1がたつことになる。出力ビットCが1の時は、出力ビットC'は0となる。出力ビットCが0の時、出力ビットC'は1となる。

【0094】各タグのタグ値は、宛先ビットの値を逆転

させた値でOR演算される。これは、処理対象となっていないタグレジスタが小さい値を持っている場合に、現在処理中の宛先に関するタグレジスタの比較結果に、悪影響を及ぼすことを防ぐことを目的としている。宛先ビットが0の場合は、逆転させた値が1となり、この1とタグレジスタの各ビットのOR演算すると、全ビットが1となり、最大のタグ値を示すことになる。このようにして、宛先ビットがのエンタリーのタグ値を無視することができる。木構造のそれぞれのレベルにおいて、比較器からの出力kビットは、図面では比較器の右側に示されている。次の上位レベルの比較器の一方に入力される。比較器から出力されたC及びC'は、直前の下位レベルで選択された出力ビットとAND演算される。ここでC'は、Cの値を逆転させた値を示す。その結果、検索の結果を示す出力ビットが最上位レベルにおいて出力される。最終的に1にセットされた出力ビットは、1つだけになる。即ち、最小値のタグを持つキューエンタリーに対応するビットであり、かつ、キューの先頭に最も近いキューエンタリーに対応するビットである。宛先ビットが1つも一致しない場合には、出力ビットが1つもセットされない。

【0095】Nをキューエンタリー総数とすると、この二分木からなる比較器の最長パスは、 $k * \log_2 N$ となる。前述したリニアサーチ回路と同様の考え方を利用すれば、16ビットのタグを有する128個のキューエンタリーから結果を選択するまでにかかる合計時間は、 $5.5 * \log_2 128 = 5.5 * 7 = 38.5$ ナノセカンドとなる。この時間は、キューのデータ数が2倍になっても、僅か5.5ナノセカンド延長されるだけである。例えば、256個のキューエンタリーから結果を選択するまでにかかる合計時間は、 $5.5 * \log_2 256 = 5.5 * 8 = 44.0$ ナノセカンドとなる。

【0096】図8に示した回路においては、 $N * \log_2 N$ 個の論理積ゲートが、最小値のタグを持った出力を選択するのに必要であった。この数は、各比較器からの出力C及びC'を次の下位層レベルにフィードバックすることで、 $2 * N$ 個まで減少させることが可能である。図9は、2つの層からなる変形例を示す図である。図9において、使用されている符号で図8と同一符号のものは相当部分である。今まで示してきた図の中では選択出力が右向きであったのに対して、この図においては左向きとなっている。この回路においては、図8に示した回路よりもより多くの時間が必要となる。なぜならば、C及びC'が選択出力に加えられなければならないからである。この余分にかかる時間は、 $\log_2 N$ 個の論理積回路の時間である。そして、論理積ゲート98の出力は、論理積ゲート102に接続され、論理積ゲート100の出力は、論理積ゲート104にそれぞれ接続される。これらの論理積ゲートには、図に示すように、宛先ビット106及び108もまた入力される。

【0097】上で述べた二分木検索方法は、スピードが速く、また、リニアサーチ方式と同じ数の比較器を使用するものである。次に、より処理速度が速い変形例について述べる。図10は、前述した2つの方式を組み合わせた構成図である。図において、キューは、L個のエントリーからなる複数のグループに分割されている。それぞれのグループには、図4に示した形式のリニアサーチ回路120が備えられている。これらのグループは、並列に検索される。その結果、グループ毎にグループの中で最小値のタグを持つエントリーを選択する。各グループからの出力124は、追加のリニアサーチ回路126に接続され、そこで全体の最小値のタグを持つグループが選択される。リニアサーチ回路126からの出力は、ビットの集まりからなる。そのビットの集まりの中の1ビットだけが、1の値を持っている。即ち、そのビットは、全体の中で最小のタグを持つグループに対応している。これらのビットの集まりは、次に論理積ゲート128の配列に接続される。この論理積ゲート128は、リニアサーチ回路120に対応している。リニアサーチ回路120の各出力ビットは、リニアサーチ回路126の出力ビットとAND演算される。その結果、最終的に最小値のタグを持ち、かつ、処理対象の宛先ビットに1がたっているキューエントリーが選択される。

【0098】以下に、より詳細に説明する。この結合方式においては、L個のエントリーからなるグループが並列して比較され、並列してリニアサーチ方式で比較される。そして、その結果が上位層のレベルで更に比較される。従って、各グループ毎にそのグループ内の最小値のタグを持つkビットの値を生成する。更に、L本の出力ラインを生成する。木構造の中の末端でないグループの出力ラインは、折り返され、次の下位層レベルのグループの出力ラインとAND演算される。それ故、葉に当たるグループの出力ラインは、そのタグがそのグループ内でも最小であり、更に全ての上位層レベルのグループの中で最小である場合のみ1がセットされる。

【0099】各グループの最長パスの長さは、 $k+L$ である。キューエントリー総数をNとすると、グループ数は、 N/L となり、 $\log_2 N$ 個の階層数が存在する。この場合、回路全体が最終結果を算出するまでの時間は、{1階層が結果を算出する時間+出力時間}×階層数となり、 $(k+L+1) * \log_2 N$ に比例する。例えば、16ビット($k=16$)の比較器が5.5ナノ秒必要とする場合、16個のエントリー($L=16$)を持つ1つのグループは、およそkビットの走査比較時間=5.5+5.5=2*5.5=11ナノ秒を検索に要することになる。16²個=256個($N=256$, $L=16$, 16グループ)のエントリーのキューは、 $\log_2 N = \log_2 256 = 8$ となり、この数字(11ナノ秒)の倍だけの時間を必要とする。これは、約22ナノ秒である。また、木構造

の中で用いられる比較器の総数は、 2^k である。Lの値に関わらず、比較器の総数は、およそ 2^k となる。

【0100】より早い処理スピードを達成するために、必要であれば、木構造をパイプライン化することも可能である。だが、その場合には、回路もより複雑化するという欠点も発生する。階層化された比較器の各階層レベルは、パイプラインにおけるステージに相当する。第1サイクルで第1の出力リンクに対応する一組の宛先ビットが選択され、第1階層レベルの全ての比較器から第1階層レベルの結果が形成される。形成された結果は、次の第2階層レベルに出力される。そして、その次の第2サイクルにおいて、第2の出力リンクに対応した新たな一組の宛先ビットが第1階層レベルの比較器で選択される。一方、第2階層レベルの比較器では、第2サイクルにおいて、第1階層レベルの比較器からの出力、即ち、直前のパイプラインステージからの出力が比較される。 $\log_2 N$ 個のサイクルが経過した後、第1の出力リンクに対応する選択出力ビットが結果として得られる。それ以降のサイクルにおいては、1サイクル毎に各出力リンクに対応する結果がその都度得られる。このようにして、1ATMセル時間よりも、はるかに短い時間で各宛先の選択出力が得られる。

【0101】スケジューリングアルゴリズムをインプリメントするという観点から見ると、このキューイングシステムにより提供されるタグベースの検索方法は、大変柔軟性に富むものである。そして、スタティックなスケジューリング方法やダイナミックなスケジューリング方法いずれにも幅広く対応して実現可能なものである。以下に具体的に述べる。

【0102】単純優先度アルゴリズム及びレートモノトニックスケジューリングアルゴリズムにおいては、タグ値はバーチャルチャネルにより静的に割り当てられる。即ち、チャネルが設定されると、その優先度も決定される。そして、そのチャネルに到着するセルは全て同一のタグ値を割り当てられる。より小さいタグは、優先度のより高い優先度に対応しているので、スケジューリングの際に優先される。タグベースの検索は、 2^k レベルの優先度、或いは、レートリゾリューションをサポートする。

【0103】ほとんどのATMスイッチは、いくつかの種類の優先度スケジューリングをサポートしている。だが、サポートされる優先度レベルの数は、非常に小さい。具体的には、およそ2レベルか4レベル位である。今日では、同一のネットワーク環境において、多種多様なタスクを柔軟にスケジューリングするために、かなり多くの優先度レベルをサポートすることが必要であることが明らかになっている。特に、レートモノトニックアルゴリズムは、反復性のリアルタイムタスクをそのタスクの反復の度数によって分類する。より高い度数を持つタスクは、より高い優先度を割り当てられる。多数の種類のタスクをサポートするために充分なりゾリュショ

31

ンを持つことが望まれていた。

【0104】バーチャルクロックアルゴリズム及びウェイトドフェアキューイングスケジューリングは、単純なダイナミック・スケジューリング・アルゴリズムである。セルがATMスイッチに到着すると、到着時間とバーチャルチャネルの状態に基づいて、タグが計算される。同一のバーチャルチャネルの各セルは、異なるタグを付加される。セルがネットワークからどのようなサービスを受けるかも、タグを計算する場合に用いられるアルゴリズム要素の1つである。

【0105】最早デッドラインファーストスケジューリングアルゴリズムは、全てのスケジューリングアルゴリズムの中で、最も一般的なものである。一連のリアルタイムタスクがいかなる方法によりスケジューリングされようとも、それらのタスクは、最早デッドラインファーストアルゴリズムにより、スケジューリングされることが可能であるということが証明されている。バーチャルクロックアルゴリズムやウェイトドフェアキューイングアルゴリズムのように、最早デッドラインファーストアルゴリズムもダイナミック・スケジューリング・アルゴリズムである。このアルゴリズムによれば、タグ値は、各セルの到着時間及びバーチャルチャネルの状態に基づいて計算される。その結果、タグは、デッドライン値を持つ。

【0106】最後に、ラウンド・ロビン・アルゴリズムは、ハードウェアで実現するよりも、ソフトウェアで実現する方がより実現しやすいと考えられている典型的なアルゴリズムである。このアルゴリズムを、この発明のタグベースの検索で実現することが可能である。それは、そのアルゴリズムを他のダイナミック・スケジューリング・アルゴリズムと同様に取り扱うことにより可能となる。そうすることによって、全てのバーチャルチャネルに対して、1つのグローバルカウンタRが保持される。グローバルカウンタRの値は、処理されたラウンドを記憶する。また、各バーチャルチャネルv毎に、そのバーチャルチャネルの最新の受信セルのラウンド番号を示す値r_vが保持される。新しいセルを受信すると、タグは以下のように設定される。

```
tagv = new_rv
= max {R, old-rv} + 1
```

【0107】次に、セルは、このタグを付加されてキューに挿入される。セルをディスパッチ（発送）する時間になると、最小値のタグが選択される。ラウンド・ロビン・チャネルの選択されたセルのタグがグローバルカウンタRの値よりもより大きい場合、グローバルカウンタRの値は、タグの値で書き換えられ更新される。

【0108】全てのダイナミック・スケジューリング・アルゴリズムにおいて、タグは任意の実時間時計により増加する要素であり、時としてタグ値の計算結果がk-1ビットを超えることがある。その場合、キューの全てのエントリーのタグ値は、kが充分な桁数を持っていて

32

ば、タグレジスタの最上位ビット（k番目）が1となる。それ以上のオーバフローを防ぐためには、単にこれらのビットの全てをリセットしてゼロにし、処理を続けられればよい。これは、ロジカルクロック又はバーチャルクロックの時間を 2^{k-1} ユニット毎に、 2^{k-1} の値だけ進めることに等しい。より詳しく述べると、キューに残っている全てのエントリーのタグ値が 2^{k-1} を超え、新しいタグの値もまた 2^{k-1} よりも大きい時、全てのキューエントリーのタグレジスタの最上位ビットをリセットすることは、キューエントリーの順番を替えずに、キューにその時点で存在するタグ値及び以降キューに入ってくるタグ値を 2^{k-1} 減少させる効果がある。

【0109】実際のネットワークにおいては、トラヒックが何らかのパラメータにより保証された配信を要求した場合には、リアルタイムスケジューリングアルゴリズムを使い、保証された配信を要求しないトラヒックについては、他の公平、かつ、最適と思われるスケジューリングアルゴリズムを使うことが必要である。これは、タグの最上位ビットを使うことで、サポート可能である。例えば、あるスイッチのバーチャルチャネルがC個のクラスに分割されているとする。また、各クラスは隣のクラスの優先度よりも、より高い優先度であるとする。この時、クラスはタグの上位log₂ Cビットにエンコードされる。最も優先度の高いクラスは、最小の値を持つ。その結果、この発明のタグベース検索システムは、常に最も高い優先度のクラスから、最小値のタグを持つキューエントリーを選択する。

【0110】以上のように、この発明のデジタル通信ネットワーク用スイッチは、汎用のキューイングシステムを利用する。このシステムは、該当するエントリーで最小のタグを持つエントリーをキューから検索するという考えに基づいている。このようなシステムでは、リアルタイムの保証を要求するアプリケーションや、オーディオデータ、ビデオデータのような連続性のあるメディアや、迅速な応答を要求するアプリケーションをサポートする幅広いクラスのスケジューリングアルゴリズムを実現することが可能である。更に、1つのスイッチで同時に多数のスケジューリングアルゴリズムを実現できる。それにより、様々な基準に従って、様々なクラスのトラヒックをサポートできる。この発明は、プログラムにより実現してもよい。つまり、各入出力ポートに処理速度の速いプログラム可能なマイクロプロセッサをそれぞれ備え、このマイクロプロセッサ上で実行されるソフトウェアで、ATMセルに割り当てられたタグの計算方法を実現してもよい。このシステムは、ATMネットワークで種々のリサーチ及びプロトタイプ環境を柔軟に実現する新しいレベルを築き上げる。また、システムをインプリメントする開発者をネットワークサービスに関する負荷から解放し、システムをインプリメントする開発者をアプリケーションの仕様実現に集中させることがで

きるというメリットも得られる。

【0111】以上のように、この実施例においては、FIFO アルゴリズム及びシンプルな優先度付け方式を含む幅広いクラスのセル・スケジューリング・アルゴリズムを実現するために、タグベースの検索を行う ATM ネットワークについて説明した。このスイッチングシステムにおいては、各セルはネットワークスイッチに到着すると、パイナリーの数値を持つタグを付加される。各セルのタグは、そのセルが伝送されるバーチャルチャネル、そのバーチャルチャネルのトラヒックのクラスのスケジューリングアルゴリズム、そのセル自身の特性等に関連する情報を持っている。そして、セル、タグ及び宛先情報が、伝送待ちのセルキューの末尾に入力される。この実施例においては、キューは、VLSI 回路で構成されている。

【0112】スイッチが、特定の宛先にセルを送送しようとする時、キューに記憶されているすべてのセルを対象に並列に検索が行われ、その宛先に一致する最小のタグ値を持つセルが選択される。選択されたセルは、キューから消去される。伝送され、キューの中に同一のタグを持つ複数のセルが存在する場合には、もっとも早くキューに到着したセルが選択される。キュー検索部は、最小のタグだけでなくセルが伝送される出力ポートも識別する。以上のように、スイッチキューに記憶された各セルには、パイナリーの数値を持つタグが付加される。それにより、キューは、ある特定の宛先を持つセルの中で、最小値のタグを持つセルを速やかに検索することが可能になる。この構成を取れば、様々な特徴を持つネットワークのトラヒックをサポートするほとんど全ての公表されているスケジューリングアルゴリズムを実現できる。また、シビアなリアルタイムの要求にも応えられ、連続性のあるメディアの伝送も可能で、迅速なレスポンスも実現できる。更に、この構成を取れば、独自のスケジューリング方式とアルゴリズムを持つ複数のトラヒックのクラスを、1つのネットワーク内でサポートできる。

【0113】

【発明の効果】以上のように、この発明によれば、柔軟性に富んだスケジューリングを行うことができる。また、いろいろな種類のスケジューリングアルゴリズムを採用することができる。

【0114】また、この発明によれば、宛先ビットのベクトルは、除去された形でセルが出力されるので、スイッチングシステム内部で用いた宛先ビットのベクトルが外部に出力されることはなく、スイッチングシステムのインタフェースは、従来のスイッチングシステムと同様であり、従来のシステムと互換性のあるシステムを提供することができる。

【0115】また、この発明によれば、各種スケジューリングアルゴリズムを採用することができ、ユーザの要

求に応じた、或いは、各システムの要求に応じたスケジューリングを行うことができる。

【0116】また、この発明によれば、タグという情報を用いることにより優先度を持った処理を行うことができる。

【0117】また、この発明によれば、レジスタの内部に宛先ビットとタグを保持しているので、レジスタを参照することにより、スケジューリングを行うことができる。

【0118】また、この発明によれば、サーチ手段がリニアにタグをサーチしていくので、単純な方法により出力すべきセルを検出することができる。

【0119】また、この発明によれば、サーチ手段が階層の木構造を用いてタグをサーチするので、高速にセルを検出することができる。

【0120】また、この発明によれば、スケジュール手段がタグ手段と共通のキューとサーチ手段を備えることにより、優先度付けを行ったスイッチング処理を行うことができる。

【0121】また、この発明によれば、リニアサーチ手段を備えているので、タグを順番に比べることにより、優先度の高いセルを検出することができる。

【0122】また、この発明によれば、ログリズミックサーチ手段を備えているので、優先度の高いセルを高速に検索することができる。

【0123】また、この発明によれば、リニアサーチとログリズミックサーチを任意に結合することができるので、柔軟な検索を行える。

【0124】また、この発明によれば、タグレジスタと比較器を備えているので、タグの値を比較器により順番に比較することにより、優先度の高いセルを出力することができる。

【0125】また、この発明によれば、比較器の出力の中に最小のタグ値を示す特定ビットを設けているので、この特定ビットを利用することにより、最小値を持つタグを検出することができる。

【0126】また、この発明によれば、比較器は、各タグに対して追加の入力ビットを備えているので、この追加の入力ビットにより、比較対象となるタグであるかどうかを判定することができる。

【0127】また、この発明によれば、追加の入力ビットを備えているので、この追加の入力ビットに宛先ビットの情報を与えることにより、対応するタグが比較対象となるタグであるかどうかを判定することができる。

【0128】また、この発明によれば、最小のタグ値を持つセルが2以上存在する場合でも、先着したセルを先に出力する。従って、同一の優先度を持つセルが存在する場合には、セルの逆転現象は生じない。

【0129】また、この発明によれば、階層の木構造を用いてサーチするので、高速サーチを行うことができ

る。

【0130】また、この発明によれば、比較器からの比較ビットと宛先ビットをANDゲートによりチェックしているため、最終的にただ1つのANDゲートが有意な出力を有するように構成できる。このように、比較ビットとANDゲートの採用により、簡単な回路により、優先度を持ったスケジューリングを行える。

【0131】また、この発明によれば、ANDゲートにより宛先ビットの検査と最小のタグ値があるかどうかの検査が行われるので、最終的に宛先ビットがオンであり、最小のタグ値を持つセルを検出することができる。

【0132】また、この発明によれば、2以上のセルが最小のタグ値を有する場合に、先着のセルを選択するので、同じ優先度を持つセルの逆転が生じない。

【0133】また、この発明によれば、階層的木構造を持ったサーチを行う場合に、キューの末尾よりもキューの先頭方向にあるものを優先して出力するので、先入れ、先出し方式を保ちながら、優先度を持ったスケジューリングを行うことができる。

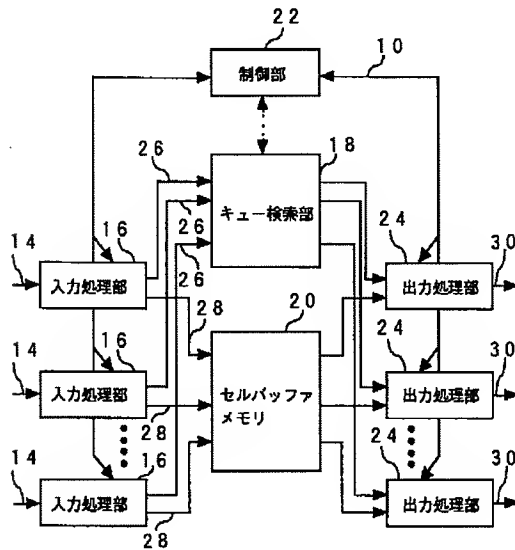
【0134】また、この発明によれば、ANDゲートの数を減少させた回路を用いてスケジューリング処理を行うことができる。

【図面の簡単な説明】

【図1】 この発明のATMスイッチの論理的構造を示すブロック図である。

【図2】 従来例のキュー検索システムを示すブロック図である。

【図1】



【図3】 図2に示す従来例の遅延優先度フィルタを備えた宛先ビットと検索回路の詳細図である。

【図4】 この発明の実施例の計算されたタグ値を選択するリニアサーチ手段に4個のキューエントリーを描いた概要図である。

【図5】 この発明のリニアサーチ手段の詳細図である。

【図6】 この発明のリニアサーチ手段の一要素のブロック図である。

【図7】 この発明の実施例の比較器の論理値の定義図である。

【図8】 この発明の実施例のログリズミックサーチ手段に8個のキューエントリーを描いた概要図である。

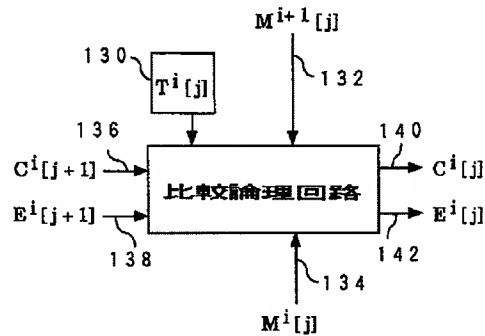
【図9】 この発明の実施例のログリズミックサーチ手段の概要図である。

【図10】 この発明の実施例のリニアサーチ手段とログリズミックサーチ手段を結合させたブロック図である。

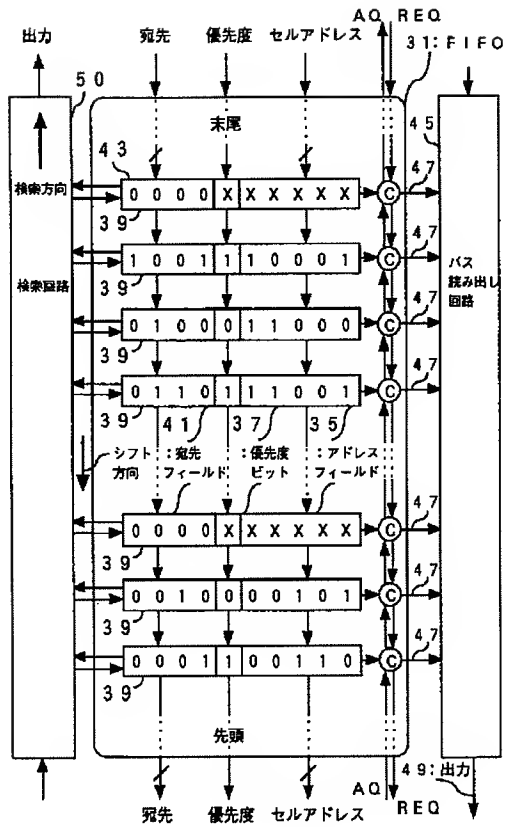
【符号の説明】

10 ATMスイッチ、14 入力リンク、16 入力処理部、18 キュー検索部、20 セルバッファメモリ、22 制御部、24 出力処理部、35 アドレスフィールド、37 優先度ビット、41 宛先フィールド、42 宛先ビット、50a 比較回路、56 出力ビット、72, 82, 84 タグレジスタ、74, 86, 88, 90, 92 比較器。

【図6】



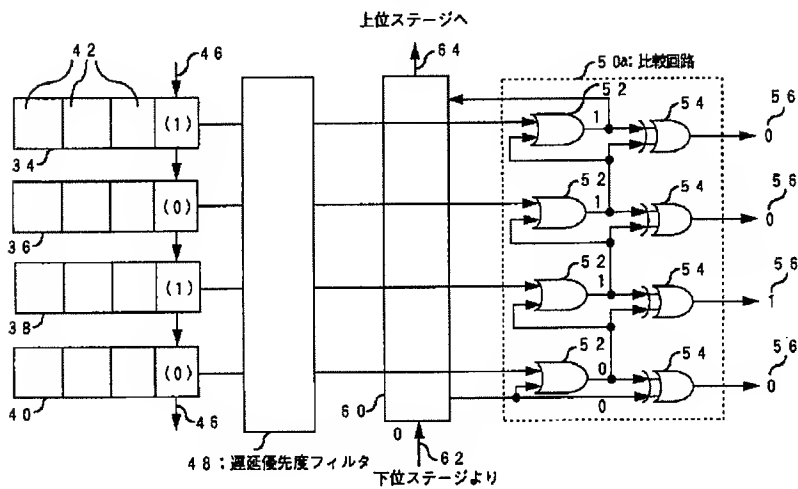
【図2】



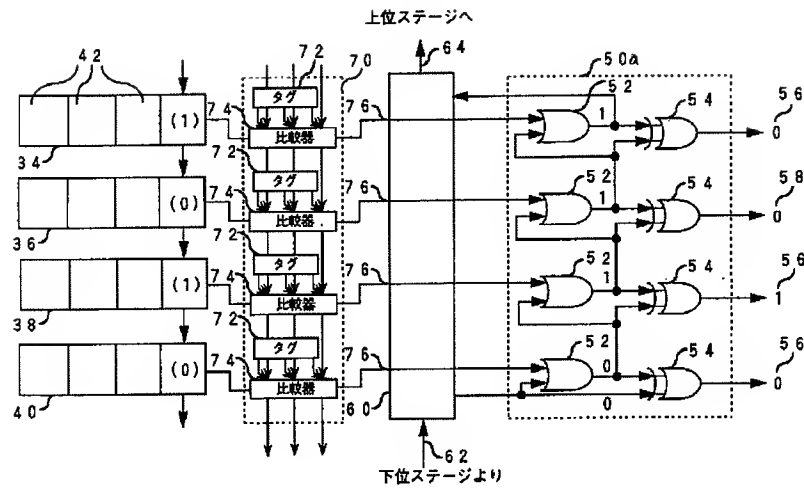
【図7】

	$C^i[j+1]$	$E^i[j+1]$	$T^i[j]$	$M^{i+1}[j]$	$M^i[j]$	$C^i[j]$	$E^i[j]$
1	0	0	0	0	0	0	0
2	0	0	0	1	1	0	0
3	0	0	1	0	0	0	0
4	0	0	1	1	1	0	0
5	0	1	X	X	X	X	X
6	1	0	0	0	0	1	0
7	1	0	0	1	0	1	0
8	1	0	1	0	1	1	0
9	1	0	1	1	1	1	0
10	1	1	0	0	0	1	1
11	1	1	0	1	0	1	0
12	1	1	1	0	0	0	0
13	1	1	1	1	1	1	1

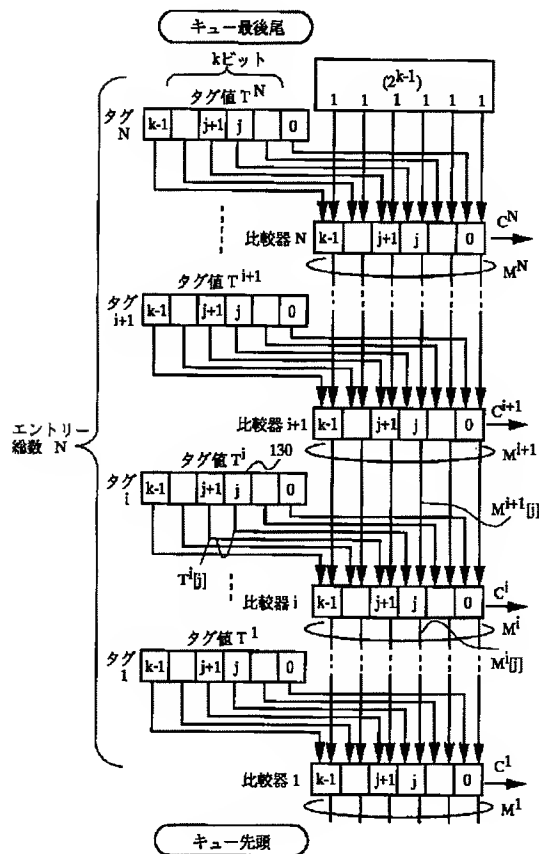
【図3】



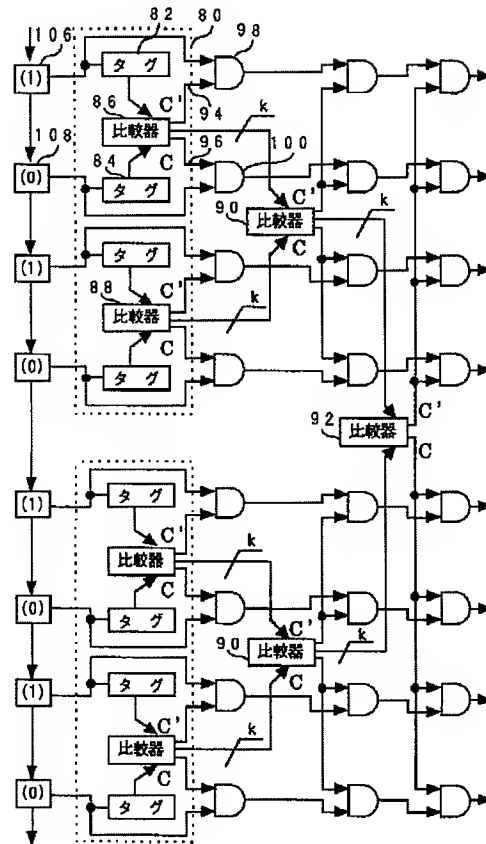
【図 4】



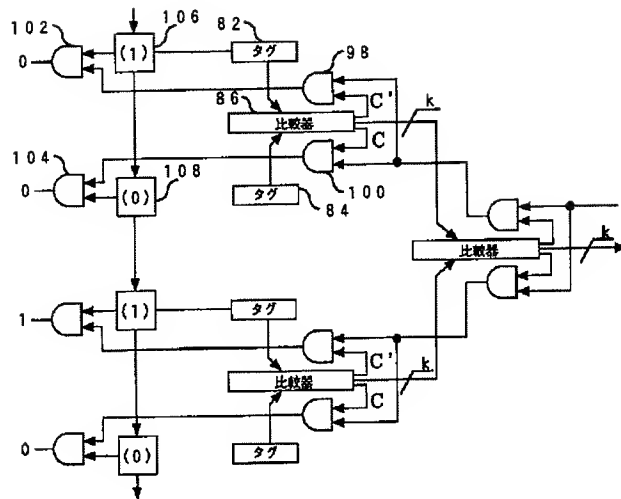
【図 5】



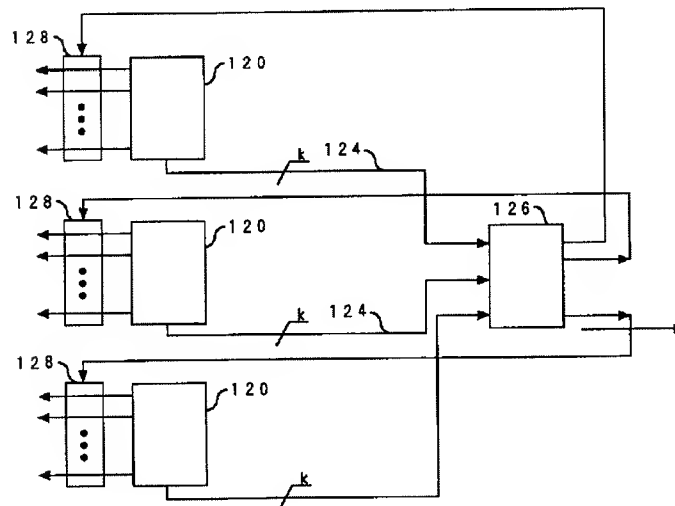
【図 8】



【図9】



【図10】



フロントページの続き

(51) Int. Cl.⁶
H04Q 3/52

識別記号 序内整理番号
101 Z 9566-5G
9466-5K

F I

技術表示箇所

H04L 11/20

102 Z